

Dynamic Programming for Markov Decision Process II

Dynamic Programming for Markov Decision Process II

Generalization

Stochastic Policy

Why stochastic policy?

Exploitation vs. Exploration

Stochastic Rewards

Markov Decision Process

Value functions

Bellman Equations

Bellman Equations for State Value

Vector-Form Bellman Equation

Bellman Equations for Q-Function

Policy Evaluation

Bellman Optimality Equation

Optimal Policies and Optimal Value Functions

Bellman Optimality Equation for State Value

Appendix

Law of total probability

Law of total expectation

Trivial Bounds of Expectation

Equivalent Formulation of Bellman equation

Notation:

- upercase (e.g. R_t): random variable
- lowercase (e.g. s_t): instance of random variable
- bold straight (e.g. \mathbf{v}): deterministic vector

Preliminary: Dynamic Programming for Markov Decision Process I

Generalization

There are multiple ways to generalize the MDP:

- deterministic policy $\pi(s) \rightarrow$ stochastic policy $\pi(a | s)$.
- deterministic rewards $r(s, a) \rightarrow$ stochastic rewards $p(r | s, a)$.
- infinite time horizon and time invariance \rightarrow finite time horizon and time dependence.

In the following, we focus on generalization to stochastic policy and stochastic rewards.

Stochastic Policy

Instead of mapping each state to a fixed action, a stochastic policy samples an action from the conditional distribution $\pi(\cdot | s)$.

$$A \sim \pi(a | s) \tag{1}$$

Remarks:

- Stochastic policy satisfies the law of total probability.

$$\sum_{a \in \mathcal{A}} \pi(a | s) = 1, \forall s \in \mathcal{S} \text{ (for finite } \mathcal{A}) \quad (2)$$

$$\int_{a \in \mathcal{A}} \pi(a | s) da = 1, \forall s \in \mathcal{S} \text{ (for continuous } \mathcal{A}) \quad (3)$$

- A deterministic policy $s \mapsto a_s$ can be seen as a special of stochastic policy

$$\pi(a | s) = \delta(a - a_s), \quad \forall a \in \mathcal{A} \quad (4)$$

where $\delta(\cdot)$ is Kronecker delta (for finite \mathcal{A}) or Dirac delta (for continuous \mathcal{A}).

- Under this setting, our stochastic policy is time-invariant (or **stationary**). i.e. The distribution $\pi(\cdot | s)$ does not depend on when we arrived at s .
- If both \mathcal{A} and \mathcal{S} , we can represent a stochastic policy as a $|\mathcal{S}| \times |\mathcal{A}|$ matrix. Each entry $\pi(a | s)$ represents the probability of taking action a at state s .

Example: An extremely simplified model of stock trading. \mathcal{S} denotes the market condition. \mathcal{A} denotes your trading action.

$$\begin{aligned} \mathcal{S} &= \{\text{bullish, bearish}\} \\ \mathcal{A} &= \{\text{buy, hold, sell}\} \end{aligned}$$

A trading strategy inspired by the principle "*Be fearful when others are greedy and to be greedy only when others are fearful*" can be modelled as a stochastic policy:

Market Condition	buy	hold	sell
bullish	0.2	0.4	0.4
bearish	0.6	0.3	0.1

Why stochastic policy?

If the agent has perfect knowledge of the environment (state transition, rewards, etc.), then any deviation from the optimal policy (which is deterministic, as shown later) will indeed yield a lower total reward. In this case, a deterministic policy suffices. The agent only needs to compute the optimal policy -- this is called **planning problem**. e.g. shortest path problem in a graph with known weights.

In contrast, a **learning problem** assumes that the agent does not have the perfect environment model. e.g. shortest path problem in a graph with unknown edge weights. The optimal policy derived from an imperfect environment model is unlikely to be optimal in real environment. In this case, stochastic policy allows the agent to explore the environment by occasionally trying different actions, even at a risk of receiving a lower total reward.

Exploitation vs. Exploration

- **Exploitation:** execute the optimal policy based on current understanding of the environment. It yields the highest total reward based on current understanding but misses the opportunity for a potentially higher total reward.
- **Exploration:** deviate from the optimal policy by occasionally trying different actions. It may yield a lower or higher total reward compared to exploitation.

Real-life analogy: Suppose you're looking for good restaurants in your city. You've found that restaurant A is your favorite spot for lunch and restaurant B is your favorite for dinner. Exploitation means you always return to A for lunch and B for dinner — you're likely to be satisfied, but you miss the chance to discover new favorites. Exploration means occasionally trying a new place, say restaurant X. It might not be as good as A, or it could turn out to be even better.

In planning problems, only exploitation is necessary because the environment is fully known. In learning problems, however, it is crucial to balance exploitation and exploration.

Stochastic Rewards

Given a state action pair (s, a) , instead of receiving a deterministic reward, the agent receives a stochastic reward sampled from $p(\cdot \mid s, a)$.

$$R \sim p(r \mid s, a)$$

Moreover, we let $\mathcal{R} \subseteq \mathbb{R}$ denote the reward sets. Then, $p(r \mid s, a)$ is either a PMF (if \mathcal{R} is finite) or a PDF (if \mathcal{R} is continuous).

Remarks:

- Stochastic reward satisfies the law of total probability

$$\sum_{r \in \mathcal{R}} p(r \mid s, a) = 1, \forall (s, a) \in \mathcal{S} \times \mathcal{A} \quad (5)$$

$$\text{or } \int_{r \in \mathcal{R}} p(r \mid s, a) dr = 1, \forall (s, a) \in \mathcal{S} \times \mathcal{A} \quad (6)$$

- A deterministic reward $(s, a) \mapsto r_{sa}$ can be seen as a special of stochastic reward

$$p(r \mid s, a) = \delta(r - r_{sa}), \quad \forall r \in \mathcal{R} \quad (7)$$

where $\delta(\cdot)$ is Kronecker delta (for finite \mathcal{R}) or Dirac delta (for continuous \mathcal{R}).

- The stochastic reward is again stationary. i.e. The distribution of the reward does not depend on when we arrive at s or when we take action a .

Markov Decision Process

A generalized Markov decision process (MDP) consists of

- \mathcal{S} : set of states.
- \mathcal{A} : set of actions.
- \mathcal{R} : set of rewards.
- $p(s' \mid s, a)$: state transition probability.
- $p(r \mid s, a)$: reward probability.
- $\gamma \in [0, 1)$: discount factor.

The agent executes a stochastic policy π . Then, we are interested in

1. policy evaluation
2. computing optimal policy

Consider the execution of π from a fixed initial state s .

$$s \xrightarrow[R]{A} s' \xrightarrow[R']{A'} \dots$$

Key observation:

- The immediate reward R is now stochastic. The stochasticity of immediate reward comes from
 1. stochasticity in action A due to stochastic policy
 2. reward probability

- Starting from the next state, the stochasticity in R' comes from
 1. stochasticity in state S' due to state transition probability
 2. stochasticity in action A' due to stochastic policy
 3. reward probability

The total (discounted) reward is defined exactly the same as in MDP I:

$$G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots \quad (8)$$

$$= \sum_{k=0}^{\infty} \gamma^k R_{t+k} \quad (9)$$

Remarks:

- For each $\tau = t, t+1, \dots$, $R_\tau \sim p(r \mid S_\tau, A_\tau)$.
- Recursive structure:

$$G_t = R_t + \gamma G_{t+1} \quad (10)$$

Value functions

State value:

$$v_\pi(s) = \mathbb{E}[G_t \mid S_t = s], \quad \forall s \in \mathcal{S} \quad (11)$$

Q-function:

$$q_\pi(s, a) = \mathbb{E}[G_t \mid S_t = s, A_t = a], \quad \forall s \in \mathcal{S}, \forall a \in \mathcal{A} \quad (12)$$

Remarks:

- $v_\pi(s) \triangleq$ the expected total reward of following a stochastic policy π from now onwards.
- $q_\pi(s, a) \triangleq$ the expected total reward of taking a **deterministic** action a now and then following a stochastic policy π .
- Relation between q_π and v_π :

$$\begin{aligned} v_\pi(s) &= \mathbb{E}_{a \sim \pi(a|s)}[q_\pi(s, a)] \\ q_\pi(s, a) &= \mathbb{E}_{r \sim p(r|s,a)}[r] + \gamma \mathbb{E}_{s' \sim p(s'|s,a)}[v_\pi(s')] \end{aligned}$$

- Proof: c.f. section *Bellman Equations for Q-Function*

Suppose the agent execute some policy π on initial state s and then follow another policy π' onwards. Then, the expected total reward is

$$\mathbb{E}[G_t \mid S_t = s] = \mathbb{E}_{a \sim \pi(a|s)}[q_{\pi'}(s, a)] \quad (13)$$

If π happens to be deterministic (i.e. $a = \pi(s)$), the total reward becomes

$$\mathbb{E}[G_t \mid S_t = s] = q_{\pi'}(s, \pi(s)) \quad (14)$$

Illustration:

$$s \xrightarrow[A]{A \sim \pi(a|s)} s' \xrightarrow[A']{A' \sim \pi'(a|s)} \dots$$

Bellman Equations

Bellman Equations for State Value

Bellman equation for v_π : For each $s \in \mathcal{S}$, it holds that

$$v_\pi(s) = \mathbb{E}_{a \sim \pi(a|s)} \left[\mathbb{E}_{r \sim p(r|s,a)}[r] \right] + \gamma \mathbb{E}_{a \sim \pi(a|s)} \left[\mathbb{E}_{s' \sim p(s'|s,a)}[v_\pi(s')] \right] \quad (\text{BE-V})$$

Equivalent formulation:

$$v_\pi(s) = \mathbb{E}_{a \sim \pi(a|s)} \left[\mathbb{E}_{r \sim p(r|s,a)}[r] + \gamma \mathbb{E}_{s' \sim p(s'|s,a)}[v_\pi(s')] \right] \quad (\text{BE-V1})$$

$$v_\pi(s) = \mathbb{E}_{r \sim p(r|s)}[r] + \gamma \mathbb{E}_{s' \sim p(s'|s)}[v_\pi(s')] \quad (\text{BE-V2})$$

For finite \mathcal{S} , \mathcal{A} and \mathcal{R} , Bellman equation becomes

$$v_\pi(s) = \sum_a \pi(a|s) \sum_r p(r|s,a) \cdot r + \gamma \sum_a \pi(a|s) \sum_{s'} p(s'|s,a) \cdot v_\pi(s') \quad (15)$$

$$= \sum_a \pi(a|s) \left[\sum_r p(r|s,a) \cdot r + \gamma \sum_{s'} p(s'|s,a) \cdot v_\pi(s') \right] \quad (16)$$

Remarks:

- Bellman equation holds for all $s \in \mathcal{S}$. For finite state space, there are $|\mathcal{S}|$ equations in total. There are further equivalent formulation of Bellman equations in the Appendix.
- In (BE-V): The 1st term represents the expected immediate reward. The 2nd term in (BE-V) represents the discounted expected future reward.
- In (BE-V2): The distributions $p(r|s)$ and $p(s'|s)$ are often denoted by $p_\pi(r|s)$ and $p_\pi(s'|s)$ since they both implicitly depend on π . Specifically, they can be obtained from system parameters by applying the law of total probability (c.f. Appendix)

$$p_\pi(r|s) = \mathbb{E}_{a \sim \pi(a|s)}[p(r|s,a)] \quad (17)$$

$$p_\pi(s'|s) = \mathbb{E}_{a \sim \pi(a|s)}[p(s'|s,a)] \quad (18)$$

- Illustration:

$$s \xrightarrow[R_t]{A_t} S_{t+1} \xrightarrow[R_{t+1}]{A_{t+1}} \dots$$

Proof. We first show the validity of (BE-V2). Then, we show (BE-V) and (BE-V1) are equivalent to (BE-V2).

Plugging $G_t = R_t + \gamma G_{t+1}$ into the value function, we get

$$\begin{aligned} v_\pi(s) &= \mathbb{E}[R_t + \gamma G_{t+1} | S_t = s] \\ &= \underbrace{\mathbb{E}[R_t | S_t = s]}_{\text{(I)}} + \underbrace{\gamma \mathbb{E}[G_{t+1} | S_t = s]}_{\text{(II)}} \end{aligned}$$

(I) \triangleq expected immediate reward:

$$\underbrace{\mathbb{E}[R_t | S_t = s]}_{\text{(I)}} = \mathbb{E}_{r \sim p(r|s)}[r]$$

(II) \triangleq expected future reward:

$$\begin{aligned}\underbrace{\mathbb{E}[G_{t+1} \mid S_t = s]}_{(II)} &= \mathbb{E}_{s' \sim p(s'|s)} \left[\mathbb{E}[G_{t+1} \mid S_t = s, S_{t+1} = s'] \right] \\ &= \mathbb{E}_{s' \sim p(s'|s)} \left[\mathbb{E}[G_{t+1} \mid S_{t+1} = s'] \right] \\ &= \mathbb{E}_{s' \sim p(s'|s)} [v_\pi(s')]\end{aligned}$$

Combining (I) and (II), we obtain (BE-V2):

$$v_\pi(s) = \underbrace{\mathbb{E}_{r \sim p(r|s)}[r]}_{(I)} + \gamma \underbrace{\mathbb{E}_{s' \sim p(s'|s)}[v_\pi(s')]}_{(II)}$$

(BE-V) and (BE-V1) follow from (BE-V2) and the law of total expectation (c.f. Appendix) as

$$\begin{aligned}\underbrace{\mathbb{E}_{r \sim p(r|s)}[r]}_{(I)} &= \mathbb{E}_{a \sim \pi(a|s)} \left[\mathbb{E}_{r \sim p(r|s,a)}[r] \right] \\ \underbrace{\mathbb{E}_{s' \sim p(s'|s)}[v_\pi(s')]}_{(II)} &= \mathbb{E}_{a \sim \pi(a|s)} \left[\mathbb{E}_{s' \sim p(s'|s,a)}[v_\pi(s')] \right]\end{aligned}$$

■

Vector-Form Bellman Equation

Recall the Bellman equation for finite \mathcal{S} , \mathcal{A} and \mathcal{R} :

$$v_\pi(s) = \sum_a \pi(a \mid s) \sum_r p(r \mid s, a) \cdot r + \gamma \sum_a \pi(a \mid s) \sum_{s'} p(s' \mid s, a) \cdot v_\pi(s')$$

Switching the orders of sums in both terms, we get

$$\begin{aligned}v_\pi(s) &= \sum_r \underbrace{\sum_a \pi(a \mid s) p(r \mid s, a)}_{p_\pi(r|s)} \cdot r + \gamma \sum_{s'} \underbrace{\sum_a \pi(a \mid s) p(s' \mid s, a)}_{p_\pi(s'|s)} \cdot v_\pi(s') \\ &= \underbrace{\sum_r p_\pi(r \mid s) \cdot r}_{r_\pi(s)} + \gamma \sum_{s'} p_\pi(s' \mid s) \cdot v_\pi(s')\end{aligned}$$

Assume $\mathcal{S} = \{\varsigma_1, \dots, \varsigma_{|\mathcal{S}|}\}$. Consider $v_\pi(\varsigma_i), i = 1, \dots, |\mathcal{S}|$ as unknowns. For a certain π , the 1st term $r_\pi(s)$ is just a known number. The 2nd term is a linear combination of $\varsigma_1, \dots, \varsigma_{|\mathcal{S}|}$. Therefore, we have

$$\underbrace{\begin{bmatrix} v_\pi(\varsigma_1) \\ v_\pi(\varsigma_2) \\ \vdots \\ v_\pi(\varsigma_n) \end{bmatrix}}_{\mathbf{v}_\pi} = \underbrace{\begin{bmatrix} r_\pi(\varsigma_1) \\ r_\pi(\varsigma_2) \\ \vdots \\ r_\pi(\varsigma_n) \end{bmatrix}}_{\mathbf{r}_\pi} + \gamma \underbrace{\begin{bmatrix} p_\pi(\varsigma_1 \mid \varsigma_1) & \dots & p_\pi(\varsigma_n \mid \varsigma_1) \\ p_\pi(\varsigma_1 \mid \varsigma_2) & \dots & p_\pi(\varsigma_n \mid \varsigma_2) \\ \vdots & \dots & \vdots \\ p_\pi(\varsigma_1 \mid \varsigma_n) & \dots & p_\pi(\varsigma_n \mid \varsigma_n) \end{bmatrix}}_{\mathbf{P}_\pi} \cdot \underbrace{\begin{bmatrix} v_\pi(\varsigma_1) \\ v_\pi(\varsigma_2) \\ \vdots \\ v_\pi(\varsigma_n) \end{bmatrix}}_{\mathbf{v}_\pi}$$

Bellman equation (vector form)

$$\mathbf{v}_\pi = \mathbf{r}_\pi + \gamma \mathbf{P}_\pi \mathbf{v}_\pi \quad (19)$$

Remarks:

- In finite settings, the Bellman equation has again the same vector form, even though policy and rewards are stochastic. Here, the \mathbf{r}_π is a vector of expected immediate rewards w.r.t. $r \sim p_\pi(r \mid s)$, which is the statistical average of $p_\pi(r \mid s, a)$ over $a \sim \pi(a \mid s)$. The row-stochastic matrix \mathbf{P}_π consists of state transition probabilities $p_\pi(s' \mid s)$ under π , which is the statistical average of $p_\pi(s' \mid s, a)$ over $a \sim \pi(a \mid s)$

- The vector-form can again be solve analytically. However, analytical solution is rarely used unless the state space is small. In practice, iterative method is more frequently used (detailed later).

Bellman Equations for Q-Function

Relation between q_π and v_π :

$$v_\pi(s) = \mathbb{E}_{a \sim \pi(a|s)}[q_\pi(s, a)] \quad (\text{Q2V})$$

$$q_\pi(s, a) = \mathbb{E}_{r \sim p(r|s, a)}[r] + \gamma \mathbb{E}_{s' \sim p(s'|s, a)}[v_\pi(s')] \quad (\text{V2Q})$$

Bellman equation for q_π :

$$q_\pi(s, a) = \mathbb{E}_{r \sim p(r|s, a)}[r] + \gamma \mathbb{E}_{s' \sim p(s'|s, a)}[\mathbb{E}_{a' \sim \pi(a'|s')}[q_\pi(s', a')]] \quad (\text{BE-Q})$$

Proof: Equation (Q2V) follows from the law of total expectation:

$$\begin{aligned} v_\pi(s) &= \mathbb{E}[G_t \mid S_t = s] \\ &= \mathbb{E}_{a \sim \pi(a|s)} \left[\underbrace{\mathbb{E}[G_t \mid S_t = s, A_t = a]}_{q_\pi(s, a)} \right] \end{aligned} \quad \blacksquare$$

To show (V2Q), we use $G_t = R_t + \gamma G_{t+1}$:

$$\begin{aligned} q_\pi(s, a) &= \mathbb{E}[R_t + \gamma G_{t+1} \mid S_t = s, A_t = a] \\ &= \underbrace{\mathbb{E}[R_t \mid S_t = s, A_t = a]}_{\text{(I)}} + \gamma \underbrace{\mathbb{E}[G_{t+1} \mid S_t = s, A_t = a]}_{\text{(II)}} \end{aligned}$$

The expected immediate reward (I) is

$$\underbrace{\mathbb{E}[R_t \mid S_t = s, A_t = a]}_{\text{(I)}} = \mathbb{E}_{r \sim p(r|s, a)}[r]$$

The expected future reward (II) is

$$\begin{aligned} \underbrace{\mathbb{E}[G_{t+1} \mid S_t = s, A_t = a]}_{\text{(II)}} &= \mathbb{E}_{s' \sim p(s'|s, a)} \left[\mathbb{E}[G_{t+1} \mid S_t = s, A_t = a, S_{t+1} = s'] \right] \\ &= \mathbb{E}_{s' \sim p(s'|s, a)} \left[\mathbb{E}[G_{t+1} \mid S_{t+1} = s'] \right] \\ &= \mathbb{E}_{s' \sim p(s'|s, a)}[v_\pi(s')] \end{aligned}$$

Combining (I) and (II), we obtain (Q2V):

$$v_\pi(s) = \underbrace{\mathbb{E}_{r \sim p(r|s, a)}[r]}_{\text{(I)}} + \gamma \underbrace{\mathbb{E}_{s' \sim p(s'|s, a)}[v_\pi(s')]}_{\text{(II)}} \quad \blacksquare$$

Plugging (Q2V) into (V2Q) yields (BE-Q). \blacksquare

Policy Evaluation

In this section, we assume that \mathcal{S} , \mathcal{A} and \mathcal{R} are all finite. Under this setting, Bellman equation can be solved (either analytically or numerically) entirely based on the system model.

In contrast, when \mathcal{S} is continuous, the Bellman equation becomes an integral equation of v_π . It is challenging to derive model-based algorithms to solve Bellman equation without additional assumptions (e.g. all random variables are Gaussian). This difficulty arises because the expectations in the Bellman equation correspond to high-dimensional integrals, which are generally intractable to compute exactly. To address this difficulty, other techniques like function approximation are used (not detailed here).

Algorithm to compute state values:

BELLMAN UPDATE (vector form)

Initialize \mathbf{v}_0 arbitrarily.

For $n = 0, 1, \dots$, run until \mathbf{v}_n converges

$$\mathbf{v}_{n+1} = \mathbf{r}_\pi + \gamma \mathbf{P}_\pi \mathbf{v}_n \quad (20)$$

Equivalent element-wise form:

BELLMAN UPDATE (element-wise form)

Init $v_0(s)$ for all $s \in \mathcal{S}$

For $n = 0, 1, \dots$, run until $v_n(s)$ converges

For each $s \in \mathcal{S}$, do

$$v_{n+1}(s) = \sum_r p_\pi(r | s) \cdot r + \gamma \sum_{s'} p_\pi(s' | s) \cdot v_n(s')$$

Remarks:

- The correctness of Bellman update here follows the same principle as in MDP I. The stochasticity of policy and rewards does not alter the contractiveness of the affine mapping $f(\mathbf{v}) = \mathbf{r}_\pi + \gamma \mathbf{P}_\pi \mathbf{v}$, even though the definition of \mathbf{r}_π and \mathbf{P}_π are slightly different from those in MDP I. A detailed proof is omitted.
- In practice, $p_\pi(r | s)$ and $p_\pi(s' | s)$ can be pre-computed by

$$p_\pi(r | s) = \sum_a \pi(a | s) \cdot p(r | s, a)$$

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

Bellman Optimality Equation

Optimal Policies and Optimal Value Functions

Optimal policy π^* :

$$\pi^* \triangleq \arg \max_{\pi} v_\pi(s) \quad (21)$$

Optimal value function v^* :

$$v^*(s) \triangleq v_{\pi^*}(s) = \max_{\pi} v_\pi(s) \quad (22)$$

Optimal Q-function q^* :

$$q^*(s, a) \triangleq q_{\pi^*}(s, a) = \max_{\pi} q_\pi(s, a) \quad (23)$$

Remark:

- The optimization variable π is a probability distribution over \mathcal{A} . Optimization problems with PDF as variables are nasty to deal with. However, we will see later that π^* is actually deterministic and thus can be written as delta function.

Equivalent formulation of v^* :

$$\begin{aligned} v^*(s) &= \max_{\pi} \mathbb{E}_{a \sim \pi(a|s)} [q_{\pi}(s, a)] \\ &= \mathbb{E}_{a \sim \pi^*(a|s)} [q^*(s, a)] \\ &= \max_{\pi} \mathbb{E}_{a \sim \pi(a|s)} [q^*(s, a)] \end{aligned} \quad (24)$$

Remarks:

- The 1st reformulation is not useful to solve v^* and π^* since the optimization variable π appears in both expectation and in Q-function.
- The 2nd reformulation is again not useful for solving v^* and π^* since it explicitly refers to π^* .
- The 3rd reformulation serves as the starting point to solve v^* and π^* . Here, the optimization variable π only appears in the expectation. The optimal Q-function inside the expectation does not require direct referencing to π^* as $q^*(s, a) = \mathbb{E}_{r \sim p(r|s, a)} [r] + \gamma \mathbb{E}_{s' \sim p(s'|s, a)} [v^*(s')]$

Proof. The 1st reformulation follows from the definition of v^* and the relation btw. v_{π} and q_{π} . The 2nd reformulation follows from the relation btw. v^* and q^* . It remains to show the 3rd reformulation.

Note that $\forall \pi, \forall s \in \mathcal{S}, \forall a \in \mathcal{A}$,

$$q_{\pi}(s, a) \leq q^*(s, a) \leq v^*(s)$$

where the 1st inequality follows from the optimality of q^* while the 2nd inequality follows from the optimality of v^* . Then:

On one hand,

$$v^*(s) = \max_{\pi} \mathbb{E}_{a \sim \pi(a|s)} [q_{\pi}(s, a)] \leq \max_{\pi} \mathbb{E}_{a \sim \pi(a|s)} [q^*(s, a)]$$

On the other hand,

$$\mathbb{E}_{a \sim \pi(a|s)} [q^*(s, a)] \leq v^*(s), \forall \pi \implies \max_{\pi} \mathbb{E}_{a \sim \pi(a|s)} [q^*(s, a)] \leq v^*(s) \quad \blacksquare$$

Bellman Optimality Equation for State Value

Bellman Optimality Equation (element-wise form)

$$v^*(s) = \max_a q^*(s, a) \quad (25)$$

$$= \max_a \left\{ \mathbb{E}_{r \sim p(r|s, a)} [r] + \gamma \mathbb{E}_{s' \sim p(s'|s, a)} [v^*(s')] \right\} \quad (26)$$

Optimal Policy:

$$\pi^*(a | s) = \delta(a - a^*) \quad \text{where } a^* = \arg \max_{a \in \mathcal{A}} q^*(s, a) \quad (27)$$

Remarks:

- BOE has again the recursive structure.

- An optimal policy is **deterministic**. It chooses the action maximizing the optimal Q-function. If there are multiple actions maximizing q^* , we may choose any of those maximizers. Hence, the optimal policy is **not unique**.
- The optimal value function v^* is however unique due to contractive mapping theorem. (Just the same as MDP I)

Proof:

Consider a non-stationary policy π' constructed as follows

1. π' **deterministically** picks $\hat{a} \triangleq \arg \max_{a \in \mathcal{A}} q^*(s, a)$ at initial state s .
2. π' thereafter follows the optimal policy π^*

The value function of π' is thus

$$v_{\pi'}(s) = q^*(s, \hat{a}) = \max_{a \in \mathcal{A}} q^*(s, a)$$

To show: $v^*(s) = \max_a q^*(s, a) \iff$ To show: $v^*(s) = v_{\pi'}(s)$.

Since π^* is an optimal policy, it is obvious that

$$v^*(s) \geq v_{\pi'}(s)$$

On the other side, we have by the trivial bound of expectation (c.f. Appendix)

$$v^*(s) = \max_{\pi} \mathbb{E}_{a \sim \pi(a|s)} [q^*(s, a)] \leq \max_a q^*(s, a) = v_{\pi'}(s)$$

with equality if

$$\pi(a | s) = \delta(a - \hat{a}) \quad \text{where } \hat{a} = \arg \max_{\tilde{a} \in \mathcal{A}} q^*(s, \tilde{a}) \quad \blacksquare$$

Appendix

Law of total probability

$$p(x) = \mathbb{E}_{\theta \sim p(\theta)} [p(x | \theta)] \quad (28)$$

$$p(x | z) = \mathbb{E}_{\theta \sim p(\theta|z)} [p(x | z, \theta)] \quad (29)$$

Alternative notation:

$$p(x) = \mathbb{E}[p(x | \Theta)] \quad (30)$$

$$p(x | z) = \mathbb{E}[p(x | z, \Theta)] \quad (31)$$

Core Intuition: computing a probability can be seen as

1. introducing an intermediate variable θ ,
2. computing the conditional probability,
3. and then averaging the conditional probability over θ .

Proof: For the sake of simplicity, we assume all random variables are continuous. Otherwise, replace the integral with sum. By the law of marginal distribution,

$$p(x) = \int_{\theta} p(x, \theta) d\theta = \int_{\theta} p(\theta) \cdot p(x | \theta) d\theta = \mathbb{E}_{\theta \sim p(\theta)} [p(x | \theta)] \quad \blacksquare$$

For $p(x | z)$, we have

$$\begin{aligned}
p(x | z) &= \frac{p(x, z)}{p(z)} && \text{def. of } p(x | z) \\
&= \frac{1}{p(z)} \int_{\theta} p(x, z, \theta) d\theta && \text{marginalization} \\
&= \frac{1}{p(z)} \int_{\theta} p(z) \cdot p(\theta | z) \cdot p(x | z, \theta) d\theta && \text{chain rule} \\
&= \int_{\theta} p(\theta | z) \cdot p(x | z, \theta) d\theta && \text{cancel out } p(z) \\
&= \mathbb{E}_{\theta \sim p(\theta|z)} [p(x | z, \theta)] && \text{def. of } \mathbb{E}_{\theta \sim p(\theta|z)}[\cdot] \quad \blacksquare
\end{aligned}$$

Law of total expectation

$$\mathbb{E}_{x \sim p(x)} [g(x)] = \mathbb{E}_{\theta \sim p(\theta)} [\mathbb{E}_{x \sim p(x|\theta)} [g(x)]] \quad (32)$$

$$\mathbb{E}_{x \sim p(x|z)} [g(x)] = \mathbb{E}_{\theta \sim p(\theta|z)} [\mathbb{E}_{x \sim p(x|z,\theta)} [g(x)]] \quad (33)$$

Alternative notation:

$$\mathbb{E}[g(X)] = \mathbb{E}[\mathbb{E}[g(X) | \Theta]] \quad (34)$$

$$\mathbb{E}[g(X) | z] = \mathbb{E}[\mathbb{E}[g(X) | z, \Theta]] \quad (35)$$

Core Intuition: computing an expectation can be seen as

1. introducing an intermediate variable θ ,
2. computing the conditional expectation,
3. and then averaging the conditional expectation over θ .

Proof: For the sake of simplicity, we assume all random variables are continuous. Otherwise, replace the integral with sum. For $\mathbb{E}_{x \sim p(x)} [g(x)]$, we have

$$\begin{aligned}
\mathbb{E}_{x \sim p(x)} [g(x)] &= \int_x g(x) \cdot p(x) dx && \text{def. of } \mathbb{E}_{x \sim p(x)}[\cdot] \\
&= \int_x g(x) \left(\int_{\theta} p(\theta) \cdot p(x | \theta) d\theta \right) dx && \text{by } p(x) = \mathbb{E}_{\theta \sim p(\theta)} [p(x | \theta)] \\
&= \int_{\theta} p(\theta) \int_x g(x) \cdot p(x | \theta) dx d\theta && \text{switch the order} \\
&= \int_{\theta} p(\theta) \mathbb{E}_{x \sim p(x|\theta)} [g(x)] d\theta && \text{def. of } \mathbb{E}_{x \sim p(x|\theta)}[\cdot] \\
&= \mathbb{E}_{\theta \sim p(\theta)} [\mathbb{E}_{x \sim p(x|\theta)} [g(x)]] && \text{def. of } \mathbb{E}_{\theta \sim p(\theta)}[\cdot] \quad \blacksquare
\end{aligned}$$

For $\mathbb{E}_{x \sim p(x|z)} [g(x)]$, we have

$$\begin{aligned}
\mathbb{E}_{x \sim p(x|z)} [g(x)] &= \int_x g(x) \cdot p(x | z) dx && \text{def. of } \mathbb{E}_{x \sim p(x|z)}[\cdot] \\
&= \int_x g(x) \left(\int_{\theta} p(\theta | z) \cdot p(x | z, \theta) d\theta \right) dx && \text{by } p(x | z) = \mathbb{E}_{\theta \sim p(\theta|z)} [p(x | z, \theta)] \\
&= \int_{\theta} p(\theta | z) \int_x g(x) \cdot p(x | z, \theta) dx d\theta && \text{switch the order} \\
&= \int_{\theta} p(\theta | z) \mathbb{E}_{x \sim p(x|z,\theta)} [g(x)] d\theta && \text{def. of } \mathbb{E}_{x \sim p(x|z,\theta)}[\cdot] \\
&= \mathbb{E}_{\theta \sim p(\theta|z)} [\mathbb{E}_{x \sim p(x|z,\theta)} [g(x)]] && \text{def. of } \mathbb{E}_{\theta \sim p(\theta|z)}[\cdot] \quad \blacksquare
\end{aligned}$$

Trivial Bounds of Expectation

Let $X : \Omega \rightarrow \mathbb{R}^n$ be a random vector with PDF p_X and $g : \mathbb{R}^n \rightarrow \mathbb{R}$ be bounded. Then

$$\mathbb{E}[g(X)] \leq \max_x g(x) \quad (36)$$

where the inequality holds iff

$$p_X(x) = \delta(x - x^*) \quad \text{where } x^* \in \arg \max_x g(x)$$

Equivalent Formulation of Bellman equation

In some literatures, the system model is given by $p(r, s' \mid s, a)$ instead of $p(r, s' \mid s, a)$ and $p(r, s' \mid s, a)$. In this setting, Bellman equation becomes

$$v_\pi(s) = \mathbb{E}_{a \sim \pi(a|s)} \left[\mathbb{E}_{r \sim p(r|s,a)} [r] + \gamma \mathbb{E}_{s' \sim p(s'|s,a)} [v_\pi(s')] \right] \quad (\text{BE-V1})$$

$$= \mathbb{E}_{a \sim \pi(a|s)} \left[\mathbb{E}_{r, s' \sim p(r, s'|s,a)} [r + \gamma v_\pi(s')] \right] \quad (\text{BE-V1}^*)$$

$$v_\pi(s) = \mathbb{E}_{r \sim p(r|s)} [r] + \gamma \mathbb{E}_{s' \sim p(s'|s)} [v_\pi(s')] \quad (\text{BE-V2})$$

$$= \mathbb{E}_{r, s' \sim p(r, s'|s)} [r + \gamma v_\pi(s')] \quad (\text{BE-V2}^*)$$

Remark:

- The probability $p(r, s' \mid s)$ in (BE-V2*) implicitly depends on π since it is obtained by marginalizing $p(r, s' \mid s, a)$ over $a \sim \pi(a \mid s)$.

$$p(r, s' \mid s) = \mathbb{E}_{a \sim \pi(a|s)} [p(r, s' \mid s, a)]$$

Proof. This is direct result of the linearity of expectation.

$$\begin{aligned} \mathbb{E}_{x \sim p(x)} [g(x)] + \mathbb{E}_{y \sim p(y)} [h(y)] &= \mathbb{E}_{x, y \sim p(x, y)} [g(x) + h(y)] \\ \mathbb{E}_{x \sim p(x|z)} [g(x)] + \mathbb{E}_{y \sim p(y|z)} [h(y)] &= \mathbb{E}_{x, y \sim p(x, y|z)} [g(x) + h(y)] \end{aligned}$$

or equivalently in simplified notation

$$\begin{aligned} \mathbb{E}[g(X)] + \mathbb{E}[h(Y)] &= \mathbb{E}[g(X) + h(Y)] \\ \mathbb{E}[g(X) \mid Z = z] + \mathbb{E}[h(Y) \mid Z = z] &= \mathbb{E}[g(X) + h(Y) \mid Z = z] \end{aligned} \quad \blacksquare$$