

$$V^\pi(s) = \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} r_t \mid s_1 = s, \pi \right] \quad s \in \mathcal{S}$$

\uparrow r.v. \uparrow realization

$$= \mathbb{E} \left[r_1 + \sum_{t=2}^{\infty} \gamma^{t-1} r_t \mid s_1 = s, \pi \right]$$

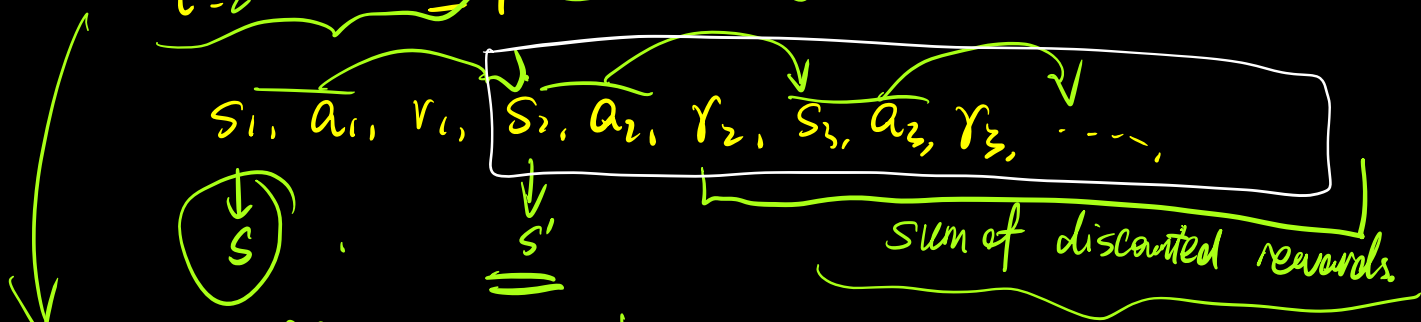
s_t : r.v.
 s, s' : realization

$$= \underbrace{\mathbb{E}[r_1 \mid s_1 = s, \pi]} + \gamma \underbrace{\mathbb{E} \left[\sum_{t=2}^{\infty} \gamma^{t-2} r_t \mid s_1 = s, \pi \right]}$$

$$= R(s, \pi(s)) + \gamma \mathbb{E}_{s_2} \left[\underbrace{\mathbb{E} \left[\sum_{t=2}^{\infty} \gamma^{t-2} r_t \mid s_1 = s, s_2, \pi \right]}_{\text{function of } s_2} \right]$$

$$= R(s, \pi(s)) + \gamma \sum_{s'} P(s' \mid s_1, \pi(s_1)) \cdot \mathbb{E} \left[\sum_{t=2}^{\infty} \gamma^{t-2} r_t \mid s_1 = s, s_2 = s', \pi \right]$$

$$\mathbb{E} \left[\sum_{t=2}^{\infty} \gamma^{t-2} r_t \mid s_1 = s, s_2 = s', \pi \right]$$



$$= \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} r_t \mid s_2 = s', \pi \right]$$

$$= \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} \cdot r_t \mid s_1 = s', \pi \right] = V^{\pi}(s').$$

Bellman Eq. for Policy Evaluation

$$\begin{aligned} V^{\pi}(s) &= R(s, \pi(s)) + \gamma \sum_{s' \in S} P(s' | s, \pi(s)) V^{\pi}(s'). \\ &= R(s, \pi(s)) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, \pi(s))} [V^{\pi}(s')]. \\ &= R(s, \pi(s)) + \gamma \langle \underbrace{P(\cdot | s, \pi(s))}_{|S| \times 1 \text{ vec.}}, \underbrace{V^{\pi}(\cdot)}_{|S| \times 1 \text{ vec.}} \rangle. \end{aligned}$$

More general: for $\pi: S \rightarrow \Delta(A)$.

replace: $R(s, \pi(s)) \rightarrow \mathbb{E}_{a \sim \pi(\cdot | s), r \sim R(s, a)} [r]$,

$s' \sim P(\cdot | s, \pi(s)) \rightarrow a \sim \pi(\cdot | s), s' \sim P(\cdot | s, a)$.

In this course: will use notation $R(s, \pi)$, $s' \sim P(\cdot | s, \pi)$ for both deterministic & stochastic policies.

$\forall s, \downarrow V^{\pi}(s) = \underbrace{R(s, \pi)} + \gamma \langle \underbrace{P(\cdot | s, \pi)}_{\Delta}, \underbrace{V^{\pi}(\cdot)}_{\text{transition matrix of MC induced by } \pi} \rangle$.

Matrix form: define.

$\rightarrow V^{\pi}$ as $|S| \times 1$ vec $[V^{\pi}(s)]_{s \in S}$.

$\rightarrow R^{\pi}$ as $[R(s, \pi)]_{s \in S}$.

$\rightarrow P^{\pi}$ as $|S| \times |S|$ matrix $[P(s' | s, \pi)]_{s \in S, s' \in S}$.

$$V^\pi = R^\pi + \gamma P^\pi V^\pi$$

unknown.

$$\underbrace{V^\pi}_{\text{unknown}} = [P(\cdot|s, \pi)]^\top \underbrace{\left[\begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \end{array} \right]}_{\Delta} \underbrace{\left[\begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \end{array} \right]}_{\Delta} = \left[\begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \end{array} \right]_s.$$

Solve: $V^\pi - \gamma P^\pi V^\pi = R^\pi$

$$(I - \gamma P^\pi) V^\pi = R^\pi$$

does inv. always exist?

$$V^\pi = \underbrace{(I - \gamma P^\pi)^{-1}}_{\Delta} R^\pi$$

YES! To prove. w.t.s. $I - \gamma P^\pi$ is non-singular.

matrix A invertible $\iff \forall x \neq \vec{0}, Ax \neq \vec{0}$

Consider arbitrary $x \neq \vec{0}$. w.t.s. $(I - \gamma P^\pi)x \neq \vec{0}$.

$$x \neq \vec{0} \iff \|x\|_\infty > 0$$

$$\|(I - \gamma P^\pi)x\|_\infty = \|x - \gamma P^\pi x\|_\infty$$

$$\geq \|x\|_\infty - \|\gamma P^\pi x\|_\infty$$

(triangular ineq. for norms).

$$= \|x\|_\infty - \gamma \|P^\pi x\|_\infty$$

$$= \|x\|_\infty - \gamma \max_s |K P(\cdot|s, \pi), x|$$

$$\geq \|x\|_\infty - \gamma \cdot \underbrace{\max_s |K P(\cdot|s, \pi)|}_{\Delta} \|x\|_\infty$$

$$= (1 - \gamma) \cdot \|x\|_\infty > 0.$$

Def. $\|\cdot\|_\infty$
 $\|f\|_\infty = \max_s |f(s)|$

apply Hölder's ineq.
 $| \langle u, v \rangle | \leq \|u\|_1 \cdot \|v\|_\infty$
 $u = P(\cdot|s, \pi)$
 $\|u\|_1 = 1$
 $v = x$
 $\|v\|_\infty = \|x\|_\infty$

Remarks: $V^\pi = (I - \gamma P^\pi)^{-1} R^\pi$ $\left\{ \begin{array}{l} R^\pi = [R(s, \pi)]_{s \in S} \\ P^\pi = [P(s' | s, \pi)]_{s, s' \in S} \end{array} \right.$

$\triangleright (1 - \gamma) \cdot (I - \gamma P^\pi)^{-1} \in \mathbb{R}^{S \times S}$.

the row vector indexed by s . is

normalized discounted occupancy of π , $(d^{\pi, s})^T \in \mathbb{R}^{1 \times |S|}$

$d^{\pi, s}(s') = (1 - \gamma) \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} \mathbb{I}[s_t = s'] \mid s_1 = s, \pi \right]$.

Equivalent def: $d_t^{\pi, s}(s') = \mathbb{E}[\mathbb{I}[s_t = s'] \mid s_1 = s, \pi]$.

marginal dist. of s_t , under π , starting from s . $= \mathbb{P}[s_t = s' \mid s_1 = s, \pi]$.

$d^{\pi, s}(s') = (1 - \gamma) \sum_{t=1}^{\infty} \gamma^{t-1} d_t^{\pi, s}(s')$.

$(I - \gamma P^\pi)^{-1} = [d^{\pi, s}]_{s \in S}^T$.

$\triangleright V^\pi = (I - \gamma P^\pi)^{-1} R^\pi$

$V^\pi(s) = \underbrace{e_s^T}_{[0, 0, \dots, 0, 1, 0, \dots, 0]^T, \uparrow s\text{-th element}} \cdot (I - \gamma P^\pi)^{-1} R^\pi = \frac{1}{1 - \gamma} \cdot (d_{\Delta}^{\pi, s})^T R^\pi$
 $= \frac{1}{1 - \gamma} \cdot \mathbb{E}_{s \sim d^{\pi, s}} [R^\pi(s)]$
 $= \frac{1}{1 - \gamma} \cdot \mathbb{E}_{s \sim d^{\pi, s}} [R(s, \pi)]$

Policy Optimization

Thm. (see Puterman '94). For infinite-horizon discounted MDPs, there always exists a stationary & deterministic policy that is optimal for all starting states simultaneously. denote such a policy as π^* .

Def. $V^* := \underline{\underline{V^{\pi^*}}}$.

Bellman Optimality Eq. $\forall s \in S$.

$$V^*(s) = \max_{a \in A} \left(R(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} [V^*(s')] \right)$$

what if we take a in s . $Q^*(s,a)$.
 $\arg\max_x = \text{optimal action} !!$

Remarks:

1. $V^\pi = R^\pi + \gamma P^\pi V^\pi$ is linear. but optimality eq. is non-linear.

2. Solution exists & is unique (will show later).
but π^* is not unique

Q: how to get π^* given V^* ??

$$\pi^*(s) = \arg\max_{a \in A} \left(R(s,a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} [V^*(s')] \right).$$

if there are ties, π^* is not unique.

can choose any action in the tie.

or even a distribution over optimal actions.

For convenience: Optimal Q-function (or action-value function)

$$Q^*(s,a) = R(s,a) + \gamma \mathbb{E}_{s' \sim p(\cdot|s,a)} [V^*(s')].$$

Why? Given Q^* , $\pi^*(s) = \underset{a}{\operatorname{argmax}} Q^*(s,a)$

Bellman Optimality Eq. (for Q^*).

$$Q^*(s,a) = R(s,a) + \gamma \mathbb{E}_{s' \sim p(\cdot|s,a)} \left[\max_{a'} Q^*(s',a') \right].$$

Remarks:

1. $V^*(s) = \max_{\pi} \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} r_t \mid s_1 = s, \pi \right]$
 $\pi \rightarrow$ unrestricted

$$Q^*(s,a) = \max_{\pi} \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} r_t \mid s_1 = s, a_1 = a, a_{2:\infty} \sim \pi \right]$$

2. Bellman opt. eq: $V^*(s) = \max_{a \in A} \left(R(s,a) + \gamma \mathbb{E}_{s' \sim p(\cdot|s,a)} [V^*(s')] \right)$

Bellman eq for PE: $V^\pi(s) = \left[R(s,\pi) + \gamma \mathbb{E}_{s' \sim p(\cdot|s,\pi)} [V^\pi(s')] \right]$

3. can also define Q^π (analogous to Q^*) evaluate at $a = \pi(s)$.

$$Q^\pi(s,a) = R(s,a) + \gamma \mathbb{E}_{s' \sim p(\cdot|s,a)} \left[Q^\pi(s', \frac{\pi}{s'}) \right].$$

* compare to Bellman eq. for Q^*

$$\left[\begin{array}{l} V^\pi(s) = \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} r_t \mid s_1 = s, \pi \right] \\ Q^\pi(s,a) = \mathbb{E} \left[\sum_{t=1}^{\infty} \gamma^{t-1} r_t \mid s_1 = s, a_1 = a, a_{2:\infty} \sim \pi \right] \end{array} \right.$$

How to solve for V^*/Q^* [Value Iteration (V2)]

VI for solving Q^* :

· Definition: Bellman optimality operator $T: \mathbb{R}^{S \times A} \rightarrow \mathbb{R}^{S \times A}$

$$\forall f \in \mathbb{R}^{S \times A}, (Tf) \in \mathbb{R}^{S \times A}$$

$$(Tf)(s, a) = R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[\max_{a'} f(s', a') \right]$$

Bellman opt. eq. for Q^* : $Q^* = TQ^*$

(Q^* is the fixed point of operator T).