# 1) State-Action Distribution
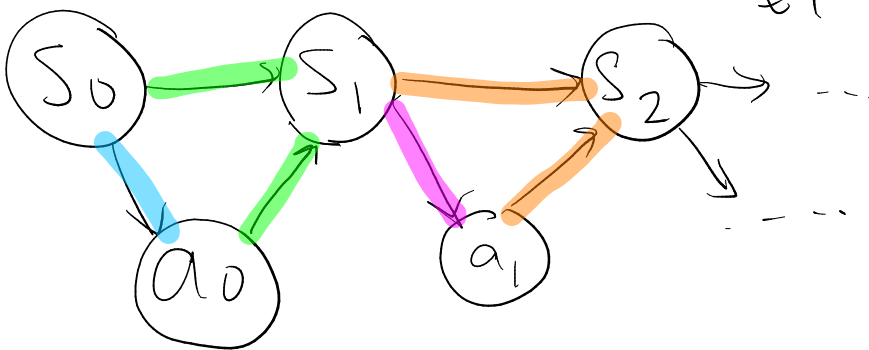
Trajectory of MDP up to step $t$:

$$(s_0, a_0, s_1, a_1, \dots, s_t, a_t)$$

What is the probability of a particular trajectory under policy $\pi$?

considering possibly stocastic policies,

$$\mathbb{P}^\pi(s_0, a_0, \dots, s_t, a_t) = \pi(a_0|s_0) P(s_1|s_0, a_0) \times$$
$$\pi(a_1|s_1) P(s_2|s_1, a_1) \times \dots$$
$$\times P(s_t|s_{t-1}, a_{t-1}) \pi(a_t|s_t)$$



← This is a graphical model of transitions which illustrates condition independe. (markov property)

What is the probability of seeing $(s, a)$ at timestep $t$, starting from $s_0$?

$$\mathbb{P}^\pi_t(s, a; s_0) = \sum_{\substack{a_{0:t-1}, \\ s_{0:t-1}}} \mathbb{P}^\pi(s_0, a_0, \dots, s_{t-1}, a_{t-1}, s_t = s, a_t = a)$$

# Discounted Average State-Action Distribution

$$d^{\pi}_{s_0}(s,a) = (1-\gamma) \sum_{t=0}^{\infty} \gamma^t \, \mathbb{P}^{\pi}_h(s,a; s_0)$$

HW0: is this a valid distribution?

$$V^{\pi}(s_0) = \frac{1}{1-\gamma} \sum_{s,a} d^{\pi}_{s_0}(s,a) \, r(s,a)?$$

## 2) Optimal Policies & Bellman Optimality

we have $A^S$ policies — which one is optimal?

$$\pi^* = \underset{\pi}{\text{argmax}} \; \mathbb{E}\left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \;\middle|\; \begin{array}{l} s_{t+1} \sim P(s_t, a_t) \\ a_t = \pi(s_t) \end{array} \right]$$
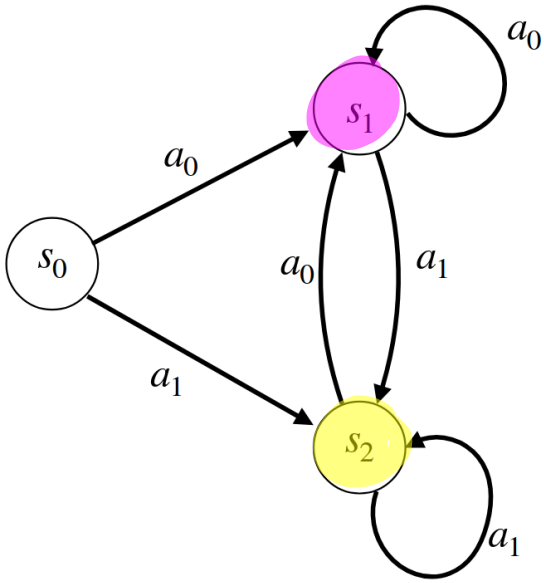
(deterministic policies & reward)

<u>Fact</u>: for infinite horizon discounted MDP, there always exists a deterministic $\pi^* : S \to A$ such that $V^{\pi^*}(s) \geq V^{\pi}(s)$ for all $s \in S$ and all $\pi$

i.e. $\pi^*$ dominates all other $\pi$ at all states! This means it is the optimal policy

notation: $V^* = V^{\pi^*}$ and $Q^* = Q^{\pi^*}$

# Example: deterministic MDP with 2 actions & 3 states

Reward is always 0 except $r(s_1, a_0) = 1$

What is the optimal policy?

consider $\pi_1(s) = a_1 \;\forall s$

$V^{\pi_1}(s_0) = V^{\pi_1}(s_1) = V^{\pi_1}(s_2) = 0$

instead, $\pi_0(s) = a_0 \;\forall s$

$V^{\pi_0}(s_0) = V^{\pi_0}(s_2) = 0 + \sum_{t=1}^{\infty} \gamma^t \cdot 1$

$= \frac{\gamma}{1-\gamma}$

$V^{\pi_0}(s_1) = \sum_{t=0}^{\infty} \gamma^t = \frac{1}{1-\gamma}$

To be rigourous we would still have to argue about the other 6 possible policies...

# Bellman Optimality

This is a key property of the optimal policy.

Theorem 1 (Bellman Optimality)

$$V^*(s) = \max_{a \in \mathcal{A}} r(s,a) + \gamma \underset{s' \sim P(s,a)}{\mathbb{E}} \left[ V^*(s') \right] \quad \text{for all } s \in \mathcal{S}$$

$Q^*(s,a)$

If we know the value of $s'$, we can use this to compute the optimal action & value of $s$.

consider a simple example with two actions & deterministic transitions

$Q^*(s,a') = r(s,a') + \gamma V(s')$

$Q^*(s,a'') = r(s,a'') + \gamma V(s'')$

$V^*(s) = \max \left\{ Q^*(s,a'), Q^*(s,a'') \right\}$

# Proof of Bellman optimality

We show for $\hat{\pi}(s) = \arg\max\limits_{a \in \mathcal{A}} Q^*(s,a)$,

that $V^{\hat{\pi}}(s) = V^*(s)$.

a) by definition of $V^*(s)$, $V^*(s) \geq V^{\hat{\pi}}(s)$ $\forall s$.

b) we now show that $V^*(s) \leq V^{\hat{\pi}}(s)$ $\forall s$.

$$V^*(s) = r(s, \pi^*(s)) + \gamma \mathop{\mathbb{E}}\limits_{s' \sim P(s, \pi^*(s))}\left[V^*(s')\right] \qquad \text{(definition of } V^*\text{)}$$

$(\text{*})$ 
$$\begin{cases} \leq \max\limits_{a} r(s,a) + \gamma \mathop{\mathbb{E}}\limits_{s' \sim P(s,a)}\left[V^*(s')\right] \qquad \left(f(a) \leq \max\limits_{u} f(u)\right) \\[2em] = r(s, \hat{\pi}(s)) + \gamma \mathop{\mathbb{E}}\limits_{s \sim P(s, \hat{\pi}(s))}\left[V^*(s')\right] \qquad \text{(definition of } \hat{\pi}\text{)} \end{cases}$$

$(\text{**})$
$$\begin{cases} = r(s, \hat{\pi}(s)) + \gamma \mathop{\mathbb{E}}\limits_{s'}\left[r(s', \pi^*(s')) + \gamma \mathop{\mathbb{E}}\limits_{s'' \sim P(s', \pi^*(s'))}\left[V^*(s'')\right]\right] \qquad \text{(defn. } V^*\text{)} \\[2em] \leq r(s, \hat{\pi}(s)) + \gamma \mathop{\mathbb{E}}\limits_{s' \sim P(s, \hat{\pi})}\left[r(s', \hat{\pi}(s')) + \gamma \mathop{\mathbb{E}}\limits_{s'' \sim P(s', \hat{\pi}(s'))}\left[V^*(s'')\right]\right] \qquad \text{(repeat } \text{*}\text{)} \end{cases}$$

$$\leq r(s, \hat{\pi}(s)) + \gamma \mathop{\mathbb{E}}\limits_{s' \sim P(s, \hat{\pi}(s))}\left[r(s', \hat{\pi}(s')) + \gamma \mathop{\mathbb{E}}\limits_{s'' \sim P(s', \pi(s'))}\left[r(s'', \hat{\pi}(s'')) + \gamma \mathop{\mathbb{E}}\limits_{s''' \sim P(s'', \hat{\pi}(s''))}\left[V^*(s''')\right]\right]\right] \qquad \text{(repeat } \text{**}\text{)}$$

$$\leq \mathop{\mathbb{E}}\limits_{s,s',s''}\left[r(s, \hat{\pi}(s)) + \gamma r(s', \hat{\pi}(s')) + \gamma^2 r(s'', \hat{\pi}(s'')) \cdots\right] \qquad \text{(repeat)}$$

$$= V^{\hat{\pi}}(s) \qquad \text{(definition } V^{\hat{\pi}}\text{)}$$

$\square$

Therefore, $V^{\hat{\Pi}}(s) = V^*(s)$ $\forall s$. This means that $\hat{\Pi}$ acheives optimal value, so

$$\hat{\Pi} = \arg\max_{a \in \mathcal{A}} Q^*(s,a) \quad \text{is an optimal policy.}$$

We now show that Bellman optimality is not only necessary, but also sufficient to characterize $V^*$.

Theorem 2:

for any $V: \mathcal{S} \rightarrow \mathbb{R}$, if $V(s) = \max_{a \in \mathcal{A}} \left[ r(s,a) + \gamma \mathop{\mathbb{E}}_{s' \sim P(s,a)} \left[ V(s') \right] \right]$ for all $s \in \mathcal{S}$, then $V(s) = V^*(s)$.

This means that finding optimal value function is equivalent to the Bellman optimality condition.

We can consider just one step between $s$ and $s'$ to check if $V = V^*$, we check if

$$\left| V(s) - \max_a \left[ r(s,a) + \gamma \mathop{\mathbb{E}}_{s' \sim P(s,a)} \left[ V(s') \right] \right] \right| = 0 \quad \forall s$$

Proof:

$$|V(s) - V^*(s)| = \left| \max_a \left[ r(s,a) + \gamma \mathop{\mathbb{E}}_{s' \sim P(s,a)}[V(s')] \right] - \max_a \left[ r(s,a) + \gamma \mathbb{E}[V^*(s')] \right] \right|$$
$$\hspace{3cm} s' \sim P(s,a)$$

basic inequalities (HW0)

$$\leq \max_a \left| r(s,a) + \gamma \mathbb{E}V(s') - r(s,a) - \gamma \mathbb{E}V^*(s') \right|$$

$$\leq \max_a \gamma \mathop{\mathbb{E}}_{s' \sim P(s,a)} \left| V(s') - V^*(s') \right|$$

$$\leq \max_a \gamma \mathop{\mathbb{E}}_{s'} \left[ \max_{a'} \mathop{\mathbb{E}}_{s'' \sim P(s',a')} |V(s'') - V^*(s')| \right]$$

(repeat)

$$\leq \max_{a_1, a_2, \ldots, a_{k-1}} \gamma^k \mathop{\mathbb{E}}_{s_k} |V(s_k) - V^*(s_k)| \longrightarrow 0 \quad \text{as } k \to \infty \quad \square$$

**Example** Recall the deterministic MDP. Now we can verify that $\pi_0(s) = a_0 \; \forall s$ is the optimal policy!

$$S_0: \left| V(S_0) - \max\left( 0 + \gamma V(S_1), \; 0 + \gamma V(S_2) \right) \right|$$
$$= \left| \frac{\gamma}{1-\gamma} - \max\left( \frac{\gamma}{1-\gamma}, \frac{\gamma^2}{1-\gamma} \right) \right| = 0$$

$$S_1: \left| V(S_1) - \max\left( 1 + \gamma V(S_1), \; 0 + \gamma V(S_2) \right) \right|$$
$$= \left| \frac{1}{1-\gamma} - \max\left( 1 + \frac{\gamma}{1-\gamma}, \frac{\gamma^2}{1-\gamma} \right) \right| = 0$$

$$S_2: \left| V(S_2) - \max\left( 0 + \gamma V(S_1), \; 0 + \gamma V(S_2) \right) \right|$$
$$= \left| \frac{\gamma}{1-\gamma} - \frac{\gamma}{1-\gamma} \right| = 0.$$

# 3) Value Iteration

How to find the optimal policy?

Algorithm: Enumeration

for all $\pi: S \to \mathcal{A}$ :

compute $V^\pi = \text{Exact-PE}(\pi)$
select $\hat{\pi}$ such that
$$V^{\hat{\pi}}(S) \geq V^\pi \; \forall \; S, \pi$$

→ a naive approach

The computation time is $O\left( A^S \cdot S^3 \right)$
Exponential complexity is a problem!

## Define Bellman operator $\mathcal{T}$:

given function $Q : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$, the Bellman operator $\mathcal{T}Q : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ defines another fn.

$$(\mathcal{T}Q)(s, a) = r(s, a) + \gamma \underset{s' \sim P(s,a)}{\mathbb{E}} \left[ \max_{a' \in \mathcal{A}} Q(s', a') \right]$$

Consider Tabular representation of $Q$,

$$Q \in \mathbb{R}^{SA}$$

$S = |\mathcal{S}|$ number of states
$A = |\mathcal{A}|$ number of actions

Then we also have

$$\mathcal{T}Q \in \mathbb{R}^{SA}$$ so we can think

of $\mathcal{T}$ as a map (nonlinear) from $\mathbb{R}^{SA}$ to $\mathbb{R}^{SA}$

## Fixed Point Motivation

By Bellman Optimality,

$$Q^*(s, a) = r(s, a) + \gamma \underset{s \sim P(s,a)}{\mathbb{E}} \max_{a'} Q^*(s', a')$$

Thus $Q^* = \mathcal{T}Q^*$ the optimal Q fn.
is a fixed point solution to $Q = \mathcal{T}Q$

## Algorithm: Value Iteration

Initialize $Q^0$

for $t = 0, 1, 2, \dots$

$$Q^{t+1} \leftarrow \mathcal{J} Q^t$$

"fixed point iteration"
like inexact Policy Iteration.

$$Q^{t+1}(s,a) \leftarrow r(s,a) + \gamma \mathop{\mathbb{E}}_{s' \sim P(s,a)} \left[ \max_{a'} Q^t(s', a') \right]$$

$\forall\ s, a$

## Convergence of Value Iteration

We will use a contraction argument.

__Lemma:__ (contraction) for any $Q, Q'$

$$\| \mathcal{J} Q - \mathcal{J} Q' \|_\infty \leq \gamma \| Q - Q' \|_\infty$$

## Proof:

$$\left| \mathcal{J} Q(s,a) - \mathcal{J} Q'(s,a) \right| = \left| \cancel{r(s,a)} + \gamma \mathop{\mathbb{E}}_{s' \sim P(s,a)} \left[ \max_{a'} Q(s', a') \right] \right.$$

$$\left. - \left( \cancel{r(s,a)} + \gamma \mathop{\mathbb{E}}_{s' \sim P(s,a)} \left[ \max_{a'} Q'(s', a') \right] \right) \right|$$

(basic inequalities HW0)

$$\leq \gamma \mathop{\mathbb{E}}_{s' \sim P(s,a)} \left| \max_{a'} Q(s', a') - \max_{a'} Q'(s', a') \right|$$

$$\leq \gamma \mathop{\mathbb{E}}_{s' \sim P(s,a)} \left[ \max_{a'} \left| Q(s', a') - Q'(s', a') \right| \right]$$

$\left( f(s') \leq \max_s f(s) \right)$

$$\leq \gamma \max_{s'} \max_{a'} \left| Q(s', a') - Q'(s', a') \right|$$

(definition of $\| \cdot \|_\infty$)

$$= \gamma \| Q - Q' \|_\infty$$

$\square$

## Lemma: (convergence) for any $Q^0$

$$\|Q^t - Q^*\|_\infty \leq \gamma^t \|Q^0 - Q^*\|_\infty$$

## Proof:

$$\|Q^t - Q^*\|_\infty = \|\mathcal{J}Q^{t-1} - \mathcal{J}Q^*\|_\infty \leq \gamma \|Q^{t-1} - Q^*\|_\infty$$
$$\leq \gamma^2 \|Q^{t-2} - Q^*\|_\infty$$
$$\cdots$$
$$\leq \gamma^t \|Q^0 - Q^*\|_\infty \quad \square$$

## From Q functions to policies

We know $\pi^*(s) = \underset{a}{\text{argmax}}\; Q^*(s,a)$

Since $Q^t(s,a) \approx Q^*(s,a)$ during value iteration,

$$\pi^t(s) = \underset{a}{\text{argmax}}\; Q^t(s,a)$$

a good choice?

---

**Theorem:** The quality of $\pi^t$ is bounded below:

$$V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty \quad \forall s \in S$$

---

## Proof:

Assume the following claim is true:

$$V^{\pi^t}(s) - V^*(s) \geq \gamma \underset{s' \sim P(s, \pi^t(s))}{\mathbb{E}}\left[V^{\pi^t}(s') - V^*(s')\right] - 2\gamma^t \|Q^0 - Q^*\|_\infty$$

Then recursing $k$ times,

$$V^{\pi^t}(s) - V^*(s) \geq \gamma^k \underset{s' \sim P(s, \pi^t(s))}{\mathbb{E}}\left[V^{\pi^t}(s') - V^*(s)\right] - 2\sum_{\ell=0}^{k} \gamma^{\ell+t} \|Q^0 - Q^*\|_\infty$$

letting $k \to \infty$,

$$V^{\pi^t}(s) - V^*(s) \geq -2\gamma^t \sum_{l=0}^{\infty} \gamma^l \|Q^0 - Q^*\|_\infty$$

$$= \frac{-2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty$$

Proof of claim:

$$V^{\pi^t}(s) - V^*(s) = Q^{\pi^t}(s, \pi^t(s)) - Q^*(s, \pi^*(s)) \quad \text{(definition)}$$
$$\underbrace{- Q^*(s, \pi^t(s)) + Q^*(s, \pi^t(s))}_{= 0}$$

$$= \gamma \underset{s' \sim P(s, \pi^t(s))}{\mathbb{E}}\left[V^{\pi^t}(s') - V^*(s')\right] + Q^*(s, \pi^t(s)) - Q^*(s, \pi^*(s))$$

$$\geq \gamma \underset{s' \sim P(s, \pi^t(s))}{\mathbb{E}}\left[V^{\pi^t}(s') - V^*(s')\right] + \underbrace{Q^*(s, \pi^t(s)) - Q^t(s, \pi^t(s))}_{\text{orange underline}} + \underbrace{\overbrace{Q^t(s, \pi^*(s)) - Q^*(s, \pi^*(s))}^{\geq 0 \text{ by } \pi^* \text{ optimality}}}_{\text{yellow underline}}$$

$$\geq \gamma \underset{s' \sim P(s, \pi^t(s))}{\mathbb{E}}\left[V^{\pi^t}(s'), V^*(s')\right] - \underbrace{\|Q^t - Q^*\|_\infty}_{} - \underbrace{\|Q^t - Q^*\|_\infty}_{}$$

by definition of $\|\cdot\|_\infty$

$$\geq \gamma \underset{s' \sim P(s, \pi^t(s))}{\mathbb{E}}\left[V^{\pi^t}(s'), V^*(s')\right] - 2\gamma^t \|Q^0 - Q^*\|_\infty \quad \text{(convergence Lemma)}$$

$\square$

# Summary of Value Iteration (VI)

1) VI (fixed point)
$$Q^{t+1} \leftarrow \mathcal{T} Q^t$$

$\xrightarrow{\text{contraction}}$

2) VI convergence
$$\|Q^t - Q^*\|_\infty \leq \gamma^t \|Q^0 - Q^*\|_\infty$$

exponentially fast
"geometric rate"

$$\pi^t(s) = \operatorname*{argmax}_a Q^t(s,a)$$

3) policy performance
$$V^{\pi^t}(s) \geq V^*(s) - \frac{2\gamma^t}{1-\gamma} \|Q^0 - Q^*\|_\infty$$

Convergence argument is similar to Iterative Policy Eval (PE)

|  |  |
|---|---|
| Bellman Eq: $$V^\pi = R + \gamma P V^\pi$$ | Bellman Optimality $$Q^* = \mathcal{T} Q^*$$ |
| Iterative PE $$V^{t+1} \leftarrow R + P V^t$$ | VI $$Q^{t+1} \leftarrow \mathcal{T} Q^t$$ |
| by contraction, $$\|V^t - V^\pi\|_\infty \leq \gamma^t \|V^0 - V^\pi\|_\infty$$ | $$\|Q^t - Q^*\|_\infty \leq \|Q^0 - Q^*\|_\infty$$ |