model-based RL: algorithms can directly call $p()$

model-free RL: only have access to $p()$ via trajectories

⌐→ { - on policy : we need trajectories using (approx)
         (SARSA)        the current agent

     - off policy : we can use any samples for
       (Q-learning)      learning

same rule for choosing
actions based on state
$\pi$, $Q$, $\varepsilon$-greedy $Q$

## Model-free methods → only access trajectories

## Temporal-difference  TD(0)  (S&B 6.1)  $\pi \rightarrow V$

$$(s_0, a_0, r_1, s_1, a_1, r_2, s_2, a_2, \ldots) \leftarrow \text{trajectory}$$

$$(s_t, a_t, r_{t+1}, s_{t+1}) \leftarrow \text{piece of trajectory}$$

$$V_\pi(s_t) = E_{a_t \sim \pi(\cdot|s_t)} E_{s_{t+1} \sim p(\cdot|s_t, a_t)} \left[ r_{t+1} + \gamma V_\pi(s_{t+1}) \right]$$

model-based method, $E \rightarrow \sum_{a_t}, \sum_{s_{t+1}}$

model-free method, sample $a_t, s_{t+1}$

~~$V_\pi(s_t) = r_{t+1} + \gamma V_\pi(s_{t+1})$  where  $(s_t, a_t, r_{t+1}, s_{t+1})$ is a sample~~

doesn't work, because one sample is not enough

$$V_\pi(s_t) \sim \underbrace{r_{t+1} + \gamma V_\pi(s_{t+1})}_{target}$$

$$V_{avg} = \frac{V_n + \overset{(1-\alpha)}{V_{n-1}} + \overset{(1-\alpha)^2}{V_{n-2}} + \dots}{\sum_i (1-\alpha)^i}$$

↰ — want to average these

— use incremental average

Very noisy improvement to $V_\pi$

$$V_\pi(s_t) = (1-\alpha) V_\pi(s_t) + \alpha \underbrace{\left( r_{t+1} + \gamma V_\pi(s_{t+1}) \right)}_{target}$$

$$\alpha = \text{learning rate} \in (0, 1)$$

$$V_\pi(s_t) = V_\pi(s_t) + \alpha \underbrace{\left( r_{t+1} + \gamma V_\pi(s_{t+1}) - V_\pi(s_t) \right)}_{\delta_t = \text{increment}}$$

where $a_t \sim \pi(\cdot | s_t)$ and $s_{t+1} \sim p(\cdot | s_t, a_t)$

↰ sampled from

# SARSA (S&B 6.4)

Use traj segments

$$S \quad A \quad R \quad S \quad A$$
$$(s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1})$$

$Q \to$ better $Q$

$\pi \not\longleftrightarrow Q$

$$\ldots, s_t, a_t, r_{t+1}, \overset{a_{t+1}}{s_{t+1} \to a_{t+1}} \to r_{t+}, s_{t+2}, a_{t+2}, r_{t+3} \ldots$$

$$Q(s_t, a_t) = r_{t+1} + \gamma E_{s_{t+1} \sim p(\cdot | s_t, a_t)} \left[ E_{a_{t+1} \sim \pi(\cdot | s_{t+1})} \left[ Q(s_{t+1}, a_{t+1}) \right] \right]$$

simplified:
$$Q(s_t, a_t) = r_{t+1} + \gamma E \left[ Q(s_{t+1}, a_{t+1}) \right]$$

$$= E \left[ \underbrace{r_{t+1} + \gamma Q(s_{t+1}, a_{t+1})}_{target} \right]$$

"on-policy" learning

$$Q(s_t, a_t) \leftarrow (1-\alpha) Q(s_t, a_t) + \alpha \left( \underbrace{r_{t+1} + \gamma Q(s_{t+1}, a'_{t+1})}_{target} \right)$$

targets should be sampled from $E$ distribution

$$a'_{t+1} = \arg\max_a Q(s_{t+1}, a)$$

$$= Q(s_t, a_t) + \alpha \underbrace{\left( r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t) \right)}_{\delta_t = increment}$$

Given $(s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1})$

Update $\boxed{Q(s_t, a_t)} \leftarrow Q(s_t, a_t) + \alpha \, \delta_t$

$(s_0, a_0, r_1, s_1, a_1)$

$(12, 3, 2, 2, 1)$

$Q = \begin{bmatrix} 5 \\ \cdot \\ 4 \end{bmatrix}$ 25

$0 \rightarrow 0.2$

$$\delta_t = r_{t+1} + \gamma \, Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)$$
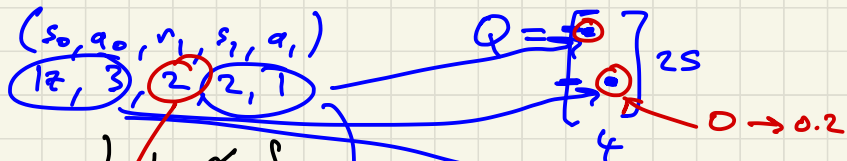
Choose actions:

greedy: $a_t = \underset{a}{\text{argmax}} \; Q(s_t, a)$  ← best choice based on current Q function (no exploration) (all exploitation)

$\varepsilon$-greedy: $a_t = \begin{cases} \text{random } a & \text{with probability } \varepsilon \\ \underset{a}{\text{argmax}} \; Q(s_t, a) & \text{otherwise} \end{cases}$

$\varepsilon$-greedy is hopefully close to optimal when $Q$ is close to optimal

exploration/ exploitation tradeoff param.

<u>Q-learning</u>   (S&B 6.5)

Use   $(s_t, a_t, r_{t+1}, s_{t+1})$

Won't assume that are chosen from $Q$

off policy

$$Q(s_t, a_t) = r_{t+1} + \gamma E_{s_{t+1} \sim p(\cdot \mid s_t, a_t)} \left[ E_{a_{t+1} \sim \pi(\cdot \mid s_{t+1})} \left[ Q(s_{t+1}, a_{t+1}) \right] \right]$$

$$a_{t+1} = \arg\max_a Q(s_{t+1}, a)$$

$$= r_{t+1} + \gamma E_{s_{t+1} \sim p(\cdot \mid s_t, a_t)} \left[ \max_a Q(s_{t+1}, a) \right]$$

Given a sample $s_{t+1}$:

$$Q(s_t, a_t) \approx r_{t+1} + \gamma \max_a Q(s_{t+1}, a)$$

$$Q(s_t, a_t) \leftarrow (1-\alpha) Q(s_t, a_t) + \alpha \left( r_{t+1} + \gamma \max_a Q(s_{t+1}, a) \right)$$

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \left( r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t) \right)$$

$$\underbrace{\qquad\qquad\qquad\qquad}_{\delta_t \text{ increment}}$$

## SARSA

$$Q(s_t, a_t) \leftarrow (1-\alpha) Q(s_t, a_t) + \alpha \left( \underbrace{r_{t+1} + \gamma Q(s_{t+1}, a_{t+1})}_{\text{target}} \right)$$

## Q-learning

$$Q(s_t, a_t) \leftarrow (1-\alpha) Q(s_t, a_t) + \alpha \left( r_{t+1} + \gamma \max_a Q(s_{t+1}, a) \right)$$

↑
don't need $a_{t+1}$

Q will improve no matter how we choose actions
(so long as we have sufficient exploration)

In practice, use $\varepsilon$-greedy actions to concentrate exploration in high-reward states.