

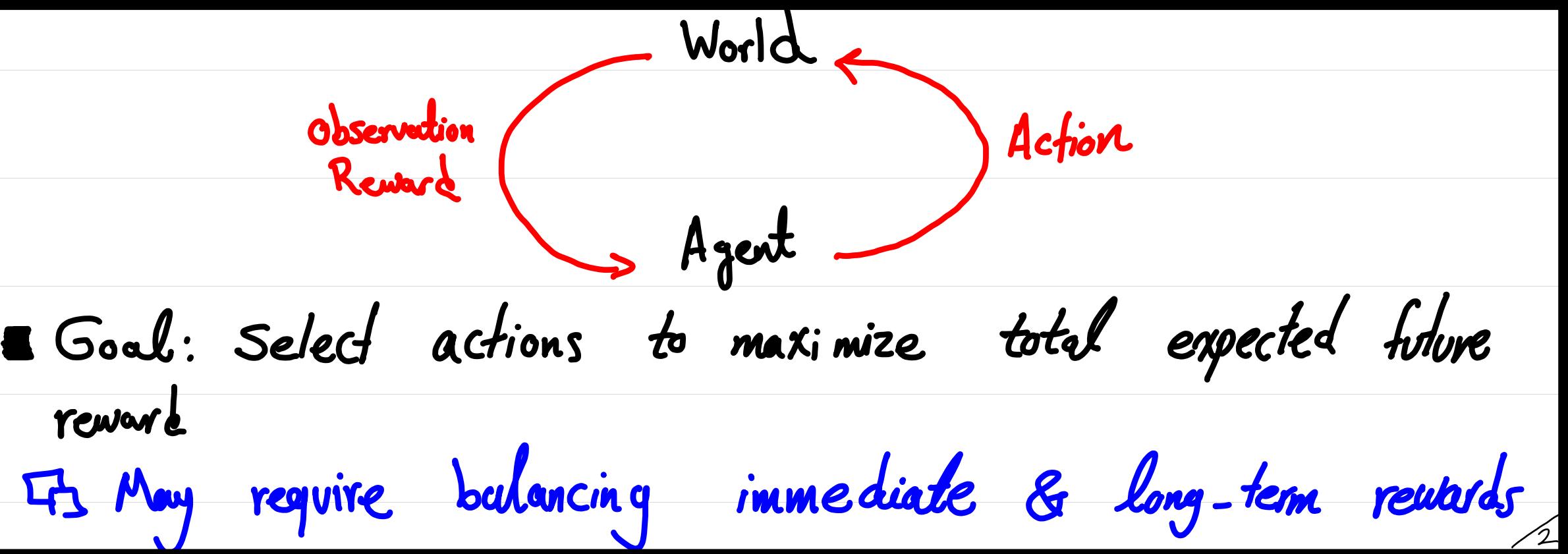
Week 11 - Part 3

■ Outcome of this part

- Return / Utility
- Policies
- Finding the "best" Policy (i.e., Solving the MDP)
 - Solving linear system of equations
 - Bellman's equations ← Next week

Returns & Episodes

■ Recall:



■ The agent's goal is to: maximize the **expected return**

↳ Return is a function of the reward sequence

$$r_0, r_1, r_2, \dots, r_T$$

Return

- In Simplest Case, the return is the sum of the rewards
 - i.e., $r_0 + r_1 + \dots + r_T$
- Such return definition makes sense in applications with finite time steps.
- However, in applications with possibly infinite time step, i.e. $T = \infty$, the above definition of return is problematic
 - as the return can be ∞ .

Return as Sum of Discounted Rewards

- That's why we use a more complex definition for return

$$U([r_0, r_1, r_2, \dots]) = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots = \sum_{k=0}^{\infty} \gamma^k r_k$$

This return is known a "Sum of discounted rewards"

$0 < \gamma < 1$ is the discount factor

If all $|r_i|$'s are bounded by α , Then $U([r_0, r_1, \dots]) \leq \alpha \frac{1}{1-\gamma}$

■ The discount parameter determines the present value of the feature rewards:

↪ A reward received K time step in feature is worth only γ^K times what would be worth if it were received immediately.

■ As γ approaches 1, the return objective takes into account the feature rewards more strongly, i.e., the agent becomes more farsighted

■ As γ approaches 0, the return objective take into account the immediate reward more strongly, i.e., the agent is more myopic (shortsighted)

Why Discounted Return is a good Model

i) As mentioned earlier, with discount factor of $\gamma < 1$, return is always bounded (i.e., finite) even when we have infinite horizon. That's because $\sum_{k=0}^{\infty} \gamma^k r_k$ is bounded if the sequence $\{r_k\}$ is finite.

Why Discounted Return is a good Model

ii) Intuitively speaking, discount return model conforms both human and animal behaviour.

We care about the reward we will get in future, but we value the immediate reward more

iii) We used return function to model agent's preference over possible outcomes, i.e., reward sequence.

A common preference assumption is stationarity, i.e.,

$$[a_1, a_2, \dots] \succ [b_1, b_2, \dots] \iff [r, a_1, a_2, \dots] \succ [r, b_1, b_2, \dots]$$

There is only one way to model a stationary preference:

$$U([r_0, r_1, \dots]) = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots, \text{ where } \gamma \in]0, 1[.$$

Policies

- A policy π , gives an action for each state
 - $\pi: S \rightarrow A$ ←
- Now that we defined return/utility as the sum of discounted rewards, we can compare policies by comparing the expected utilities obtained when they are followed
- For MDPs, we want to find an optimal policy $\pi^*: S \rightarrow A$
 - An optimal policy π^* , is the policy that maximizes the expected return if followed.
- Note: The above was the definition of deterministic Policies.
Stochastic Policies: $\pi(a|s) = P[a_t = a | s_t = s]$

Example: Erfan the Undergrad student

■ Recall the example "Erfan the undergrad".

□ $A = \{\text{"party hard"}, \text{"work hard"}, \text{"end"}\}$

$\underbrace{\text{party hard}}_{a_p}, \underbrace{\text{work hard}}_{a_w}, \underbrace{\text{end}}_{a_e}$

□ $S = \{\text{"happy"}, \text{"tired"}, \text{"burnt-out"}\}$

$\underbrace{\text{happy}}_{S_h}, \underbrace{\text{tired}}_{S_t}, \underbrace{\text{burnt-out}}_{S_b}$

■ List all possible policies for the above MDP.

□ Note that the only possible action at S_b is a_e .

$$\begin{array}{c|c|c|c} \pi_1(S_h) = a_p & \pi_2(S_h) = a_p & \pi_3(S_h) = a_w & \pi_4(S_h) = a_w \\ \pi_1(S_t) = a_p & \pi_2(S_t) = a_w & \pi_3(S_t) = a_p & \pi_4(S_t) = a_w \\ \pi_1(S_b) = a_e & \pi_2(S_b) = a_e & \pi_3(S_b) = a_e & \pi_4(S_b) = a_e \end{array}$$

Expected Return

- Following a policy yields a random episode
- The utility of a policy is the sum of discounted rewards of the episode
 - This utility is a random variable
- Let's study a simple example :

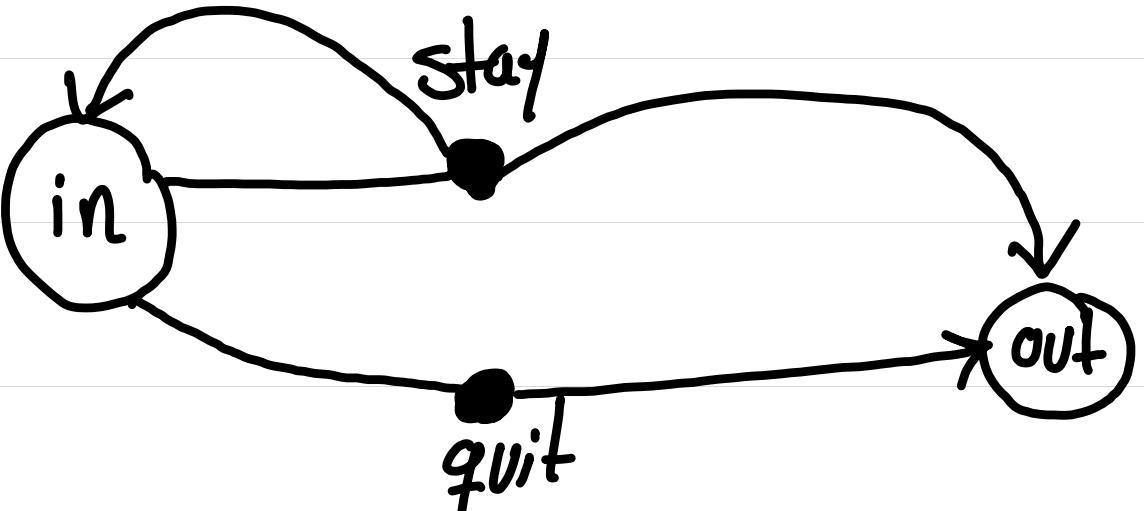
Erfan the gambler
(See the next page for the description)

Example : Erfan the Gambler

For each round $r=1, 2, \dots$

- You choose *Stay* or *quit*
- If *quit*, you get \$10 and we end the game
- If *Stay*, you get \$4 and I roll a 6-sided fair dice
 - If the dice results in 1 or 2, we end the game
 - Otherwise, Continue to the next round.

■ Let's draw the transition graph.



■ Let the discount factor be $\gamma=1$.

■ Here are some possible episodes if Erfan's policy is

to Stay:

Episode

$[(in, stay, 4), out]$

Return

4

$[(in, stay, 4), (in, stay, 4), (in, stay, 4), out]$

12

$[(in, stay, 4), (in, stay, 4), out]$

8

:

:

.

■ Let's formalize what we saw in the previous example for general MDP

■ The value (utility) of policy π at state s , denoted by $v_\pi(s)$, is the expected return if starting at s and following policy π

■ Starting at state " s " and following policy " π " can yield to different sequences of (state, action, reward), e.g.,

- $(s, \pi(s), R(s, \pi(s)))$, $(s', \pi(s'), R(s', \pi(s')))$, $(s'', \pi(s''), R(s'', \pi(s'')))$, ...
- $(s, \pi(s), R(s, \pi(s)))$, $(s', \pi(s'), R(s', \pi(s')))$, $(s'', \pi(s''), R(s'', \pi(s'')))$, ...
- $(s, \pi(s), R(s, \pi(s)))$, $(s''', \pi(s'''), R(s''', \pi(s''')))$, $(s', \pi(s'), R(s', \pi(s')))$, ...
- $(s, \pi(s), R(s, \pi(s)))$, $(s'', \pi(s''), R(s'', \pi(s'')))$, $(s', \pi(s'), R(s', \pi(s')))$, ...
- ⋮

■ Each one of the possible episodes happens with a probability

Note: in the above MDP, we assumed a reward model of type $R(s, a)$ for simplicity

* Finding the Expected Return

■ How can we formulate $V_\pi(s)$?

$$V_\pi(s) = \mathbb{E}_{(S_1, S_2, \dots) | S_0=s} [R(S, \pi(s)) + \gamma R(S_1, \pi(S_1)) + \gamma^2 R(S_2, \pi(S_2)) + \dots | (S_0=s, \pi(s))]$$

where S_1, S_2, \dots is a random sequence of states yield by following policy π . So $S_i \in \mathcal{S}$ for any $i \geq 1$.

■ What is the probability of observing the sequence S_1, S_2, S_3, \dots conditioned on $S_0 = s$? $P(S_1 | S, \pi(s)) \times P(S_2 | S_1, \pi(s)) \times \dots$

■ We can rewrite $V_\pi(s)$ as

$$V_\pi(s) = \mathbb{E}_{(S_1, S_2, \dots) | S_0=s} [R(S, \pi(s)) + \gamma (R(S_1, \pi(S_1)) + \gamma R(S_2, \pi(S_2)) + \dots) | (S_0=s, \pi(s))]$$

$$= R(s, \pi(s)) + \gamma \mathbb{E}_{(S_1, S_2, \dots) | S_0=s} [R(S_1, \pi(S_1)) + \gamma R(S_2, \pi(S_2)) + \dots | (S_0=s, \pi(s))]$$

*Finding the Expected Return

- We can rewrite $V_\pi(s)$ as

$$V_\pi(s) = R(s, \pi(s)) + \gamma \underset{(s_1, s_2, \dots) | s_0=s}{\mathbb{E}} \left[R(s_1, \pi(s_1)) + \gamma R(s_2, \pi(s_2)) + \dots \mid (s_0=s, \pi(s)) \right]$$

$\mathbb{E}[x] = \frac{1}{\gamma} \mathbb{E}_{x|y} [x|y]$

- Recall: Law of total probability $\mathbb{E}_x(X) = \sum_i P(A_i) \mathbb{E}_{x|A_i}(X|A_i)$

where $\{A_i\}$ is a finite partition of the sample space.

- Recall: Let us rewrite the above recall with an additional condition, s_0 that it is more suitable for us.

$$\mathbb{E}(X|B) = \sum_i P(A_i|B) \mathbb{E}_{x|B, A_i}(X|B, A_i)$$

- Thus, we can rewrite $V_\pi(s)$ as

$$V_\pi(s) = R(s, \pi(s)) + \gamma \sum_{s' \in S} P(s_1 = s' | s_0=s, \pi(s_0)) \underset{(s_1, s_2, s_3, \dots) | s_1=s', s_0=s}{\mathbb{E}} \left[R(s_1, \pi(s_1)) + \gamma R(s_2, \pi(s_2)) + \dots \mid (s_0=s, \pi(s)), (s_1=s', \pi(s')) \right]$$

* Finding the Expected Return

Thus, we can rewrite $V_\pi(s)$ as

$$V_\pi(s) = R(s, \pi(s)) +$$

$$\gamma \sum_{s' \in S} P(s_i = s' | s_0 = s, \pi(s_0)) \mathbb{E}_{(s_1, s_2, \dots) | s_i = s'} [R(s_1, \pi(s_1)) + \gamma R(s_2, \pi(s_2)) + \dots | (s_0 = s, \pi(s_0)), (s_i = s', \pi(s'))]$$

$$\Rightarrow V_\pi(s) = R(s, \pi(s)) +$$

$$\gamma \sum_{s' \in S} T(s, \pi(s), s') \mathbb{E}_{(s_1, s_2, \dots) | s_i = s'} [R(s', \pi(s')) + \gamma R(s_2, \pi(s_2)) + \dots | (s_i = s', \pi(s'))]$$

$$\Rightarrow V_\pi(s) = R(s, \pi(s)) + \gamma \sum_{s' \in S} T(s, \pi(s), s') V_\pi(s')$$

Finding the Expected Return

- We found a recursion for $V_\pi(s)$ under the dynamics model $T(s, a, s')$ and reward model $R(s, a)$

$$V_\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_\pi(s')$$

- More generally, we can show that for dynamics model $T(s, a, s')$ and Reward model $R(s, a, s')$

$$V_\pi(s) = \sum_{s' \in S} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_\pi(s')]$$

Review

We presented Markov Decision process and how to model it

↳ Transition Model (i.e., $T(s, a, s') \triangleq P(s'|s, a)$)

↳ Reward Model (i.e., $R(s, a, s')$)

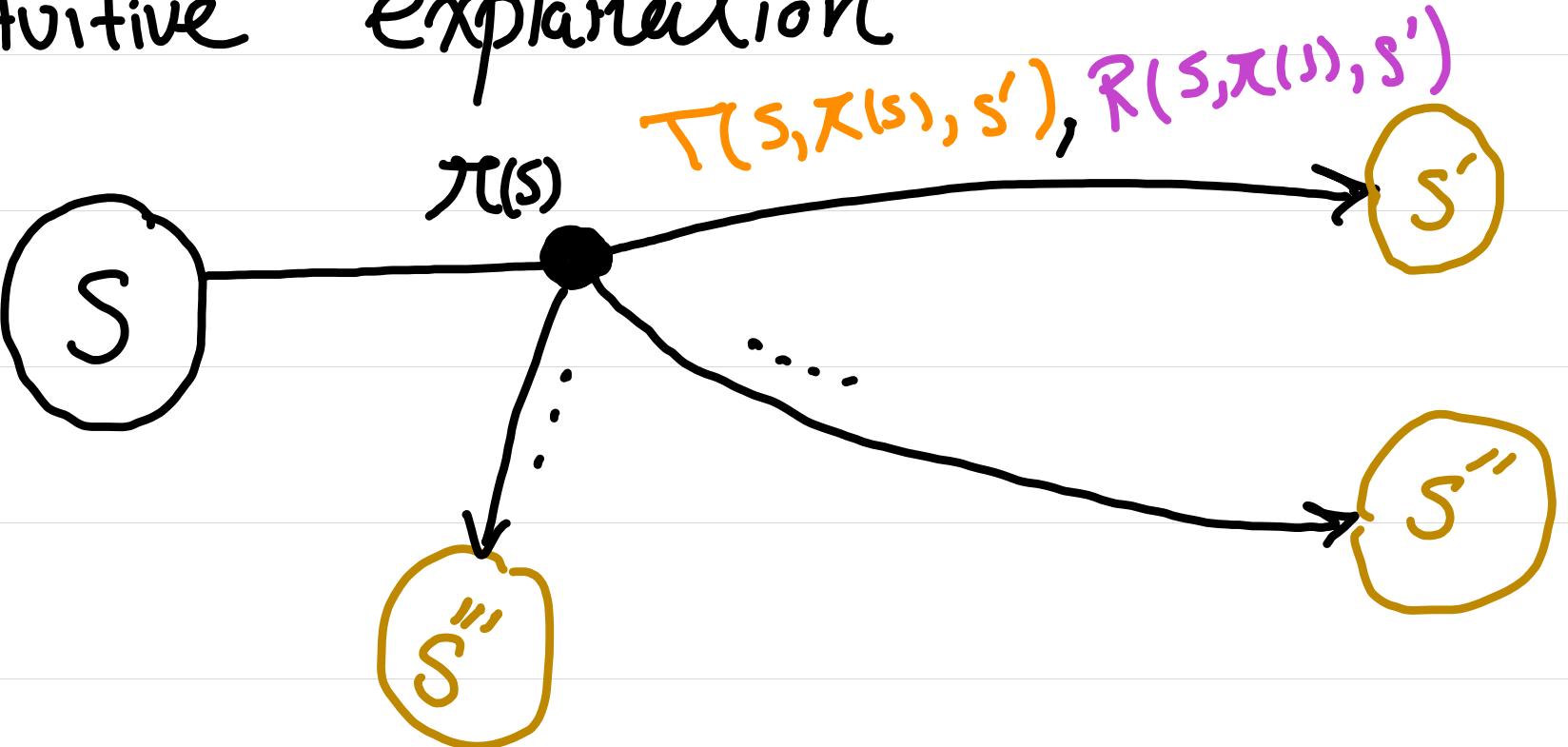
↳ γ , discount factor

↳ $V_\pi(s) \triangleq \mathbb{E}_{\substack{s_1, s_2, \dots | s}} [R(s, \pi(s), s_1) + \gamma R(s_1, \pi(s_1), s_2) + \gamma^2 R(s_2, \pi(s_2), s_3) + \dots]$

$$V_\pi(s) = \sum_{s' \in S} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_\pi(s')]$$

Finding the Expected Return

■ An intuitive explanation



with probability $T(s, \pi(s), s')$, you will receive the immediate reward $R(s, \pi(s), s')$ and go to the next state s' . From state s' , your expected utility that you would receive is $V_\pi(s')$. This is the reward you receive in the next step. Hence, you should discount it by γ .

$$V_\pi(s) = T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_\pi(s')]$$

$$+ T(s, \pi(s), s'') [R(s, \pi(s), s'') + \gamma V_\pi(s'')]$$

$$+ T(s, \pi(s), s) [R(s, \pi(s), s) + \gamma V_\pi(s)]$$

Value and Q-value of a Policy

■ We use $V_\pi(s)$ to denote the expected utility received by following policy " π " from state "S".

■ We use $Q_\pi(s, a)$ to denote the expected utility of taking action "a" from state "s", and then following policy " π ".

$$\boxed{V_\pi(s) = \sum_{s' \in S} T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_\pi(s')]}$$

↳ Value of state s in policy π

$\boxed{Q_\pi(s, a) \leftarrow \text{Q-value of state-action pair } (s, a) \text{ with policy } \pi.}$

$\boxed{\text{If you had } Q_\pi(\cdot), \text{ how would you recover } V_\pi(s)?}$

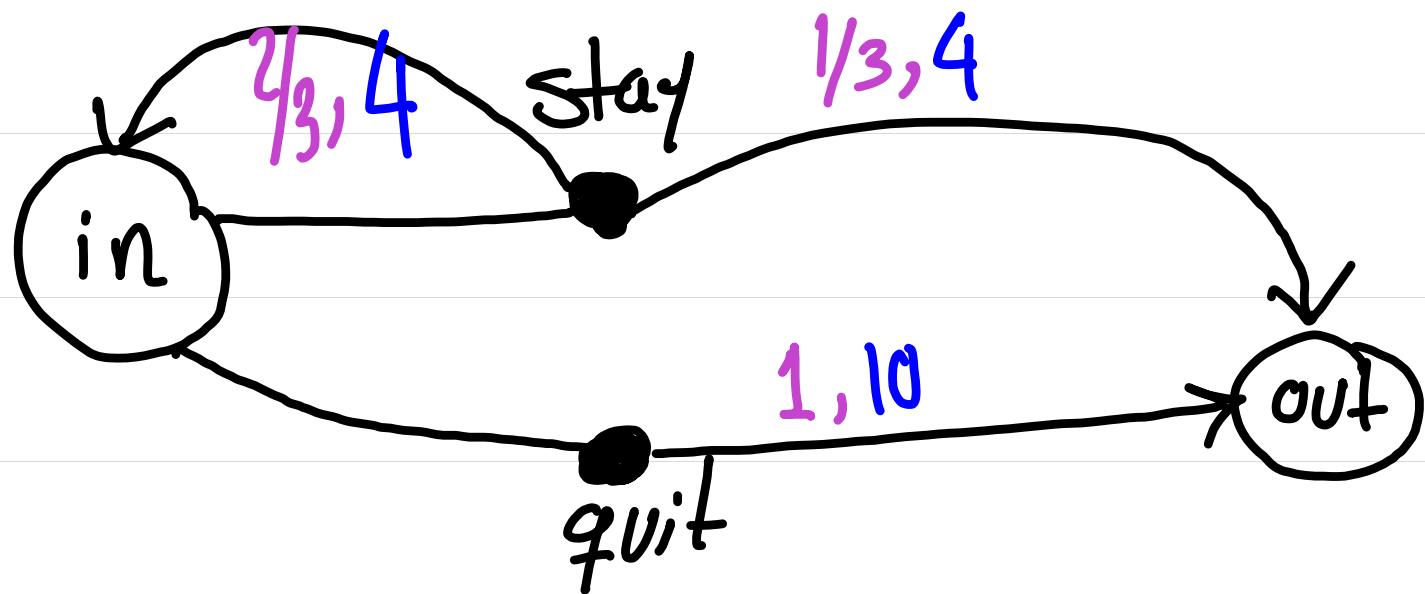
$$V_\pi(s) = Q_\pi(s, \pi(s))$$

$\boxed{\text{If you had } V(\cdot), \text{ how could you recover } Q_\pi(s, a)}$

$$Q_\pi(s, a) = \sum_{s'} T(s, a, s') [R(s, a, s') + \gamma V_\pi(s')]$$

Example: Erfan the gambler

- Let's revisit that example and find the value of "stay" policy.



- Let π_i denote the "stay" policy, i.e., $\pi_i(\text{"in"}) = \text{"stay"}$
 - Note that "out" is a terminal state and usually the value of terminal states is already known.
 - In this problem, when landing in "out" state, the game ends and we receive no more reward.
- Thus, $\forall \pi, V_\pi(\text{"out"}) = 0$, for any policy π .

■ Observe that

$$\begin{aligned}
 V_{\pi_i}("in") &= \sum_{S'} T("in", \pi_i(in), S') [R(in, \pi_i(in), S') + \gamma V_{\pi_i}(S')] \\
 &= \cancel{T(in, stay, "in")} \left[\cancel{R(in, stay, "in")} + \cancel{\gamma} V_{\pi_i}(in) \right] \\
 &\quad + \cancel{T(in, stay, out)} \left[\cancel{R(in, stay, out)} + \cancel{\gamma} \cancel{V_{\pi_i}(out)} \right]
 \end{aligned}$$

$$\Rightarrow V_{\pi_i}(in) = \frac{2}{3}[4 + V_{\pi_i}(in)] + \frac{1}{3}[4 + 0] \Rightarrow V_{\pi_i}(in) = 12.$$

$$3V_{\pi_i}(in) = 2 \times 4 + 2 \times V_{\pi_i}(in) + 1 \times 4 \Rightarrow 3V_{\pi_i}(in) - 2V_{\pi_i}(in) = 12$$

■ Let's find the value of "quit" policy.

■ Let π_2 denote the "quit" policy, i.e., $\pi_2(\text{"in"}) = \text{"quit"}$.

$$\begin{aligned} V_{\pi_2}(\text{in}) &= \sum_{s'} T(\text{in}, \pi_2(\text{in}), s') [R(\text{in}, \pi_2(\text{in}), s') + \gamma V_{\pi_2}(s')] \\ &= \cancel{T(\text{in}, \text{quit}, \text{out})} \left[10 + \cancel{V_{\pi_2}(\text{out})} \right] = 10 \end{aligned}$$

■ Which Policy do you think is better?

↪ $V_{\pi_1}(\text{in}) > V_{\pi_2}(\text{in})$. Hence, π_1 is better.

Next Lecture: Policy evaluation

■ We could find the value of the states in our previous example easily.

↳ We just had one unknown state. Easy Peasy!

■ What about larger problems? $S = \{s, s', s'', s''', \dots\}$

$$V_{\pi}(s) = T(s, \pi(s), s) [R(s, \pi(s), s) + \gamma V_{\pi}(s)] + T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_{\pi}(s')] + \dots$$

$$V_{\pi}(s') = T(s', \pi(s'), s) [R(s', \pi(s'), s) + \gamma V_{\pi}(s)] + T(s', \pi(s'), s') [R(s', \pi(s'), s') + \gamma V_{\pi}(s')] + \dots$$

$$V_{\pi}(s'') = T(s'', \pi(s'), s) [R(s'', \pi(s'), s) + \gamma V_{\pi}(s)] + T(s'', \pi(s'), s') [R(s'', \pi(s'), s') + \gamma V_{\pi}(s')] + \dots$$

Next Lecture: Policy evaluation

■ What about larger problems? $S = \{s, s', s'', s''', \dots\}$

$$V_{\pi}(s) = T(s, \pi(s), s) [R(s, \pi(s), s) + \gamma V_{\pi}(s)] + T(s, \pi(s), s') [R(s, \pi(s), s') + \gamma V_{\pi}(s')] + \dots$$

$$V_{\pi}(s') = T(s', \pi(s'), s) [R(s', \pi(s'), s) + \gamma V_{\pi}(s)] + T(s', \pi(s'), s') [R(s', \pi(s'), s') + \gamma V_{\pi}(s')] + \dots$$

$$V_{\pi}(s'') = T(s'', \pi(s'), s) [R(s'', \pi(s'), s) + \gamma V_{\pi}(s)] + T(s'', \pi(s'), s') [R(s'', \pi(s'), s') + \gamma V_{\pi}(s')] + \dots$$

\vdots

■ Let $\underline{V} = \begin{bmatrix} V_{\pi}(s) \\ V_{\pi}(s') \\ \vdots \end{bmatrix}$, $T = \begin{bmatrix} T(s, \pi(s), s) & T(s, \pi(s), s') & \dots \\ T(s', \pi(s'), s) & T(s', \pi(s'), s') & \dots \\ \vdots & \ddots & \ddots \end{bmatrix}$, $R = \begin{bmatrix} R(s, \pi(s), s) & R(s, \pi(s), s') & \dots \\ R(s', \pi(s'), s) & R(s', \pi(s'), s') & \dots \\ \vdots & \ddots & \ddots \end{bmatrix}$

■ So, $\underline{V} = (T \odot R) \underline{1} + \gamma T \underline{V}$

↳ $T \odot R$ is the hadamard product (i.e., element-wise product) of T and R .

↳ $(T \odot R) \underline{1} = \begin{bmatrix} T(s, \pi(s), s) R(s, \pi(s), s') + T(s, \pi(s), s') R(s, \pi(s), s') + \dots \\ T(s', \pi(s'), s) R(s', \pi(s'), s) + T(s', \pi(s'), s') R(s', \pi(s'), s') + \dots \\ \vdots \end{bmatrix}$

■ Thus, $\underline{V} = (I - \gamma T)^{-1} ((T \odot R) \underline{1})$

↳ Note that $(I - \gamma T)$ is invertible.