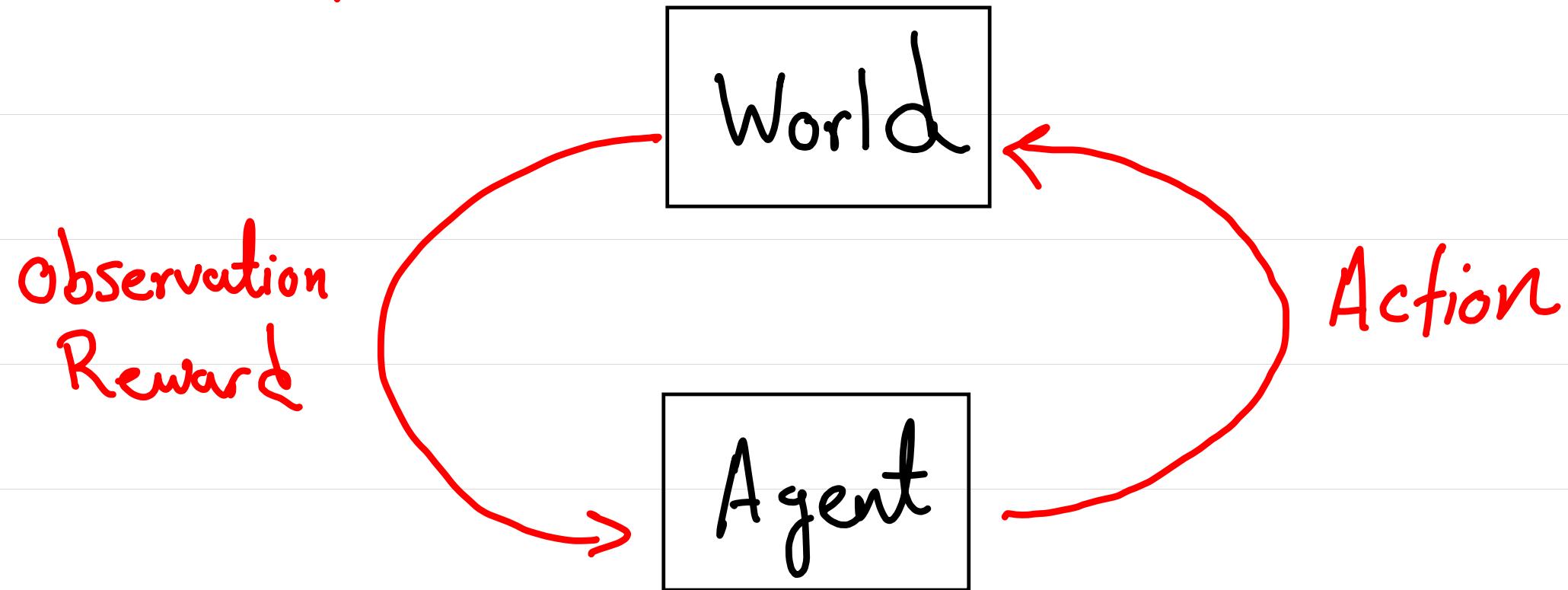


## Week 11 - Part 2

### ■ Outcome of this part

- Sequential Decision process
- Observation, history, and state
- Markov Decision Process (MDP)
  - What is Markov about MDP? Why Markov assumption is common?
- Dynamics Model & Reward Model
  - Transition Graph
- Return / Utility
- Policies
- Finding the "best" Policy (i.e., Solving the MDP) ← Next part

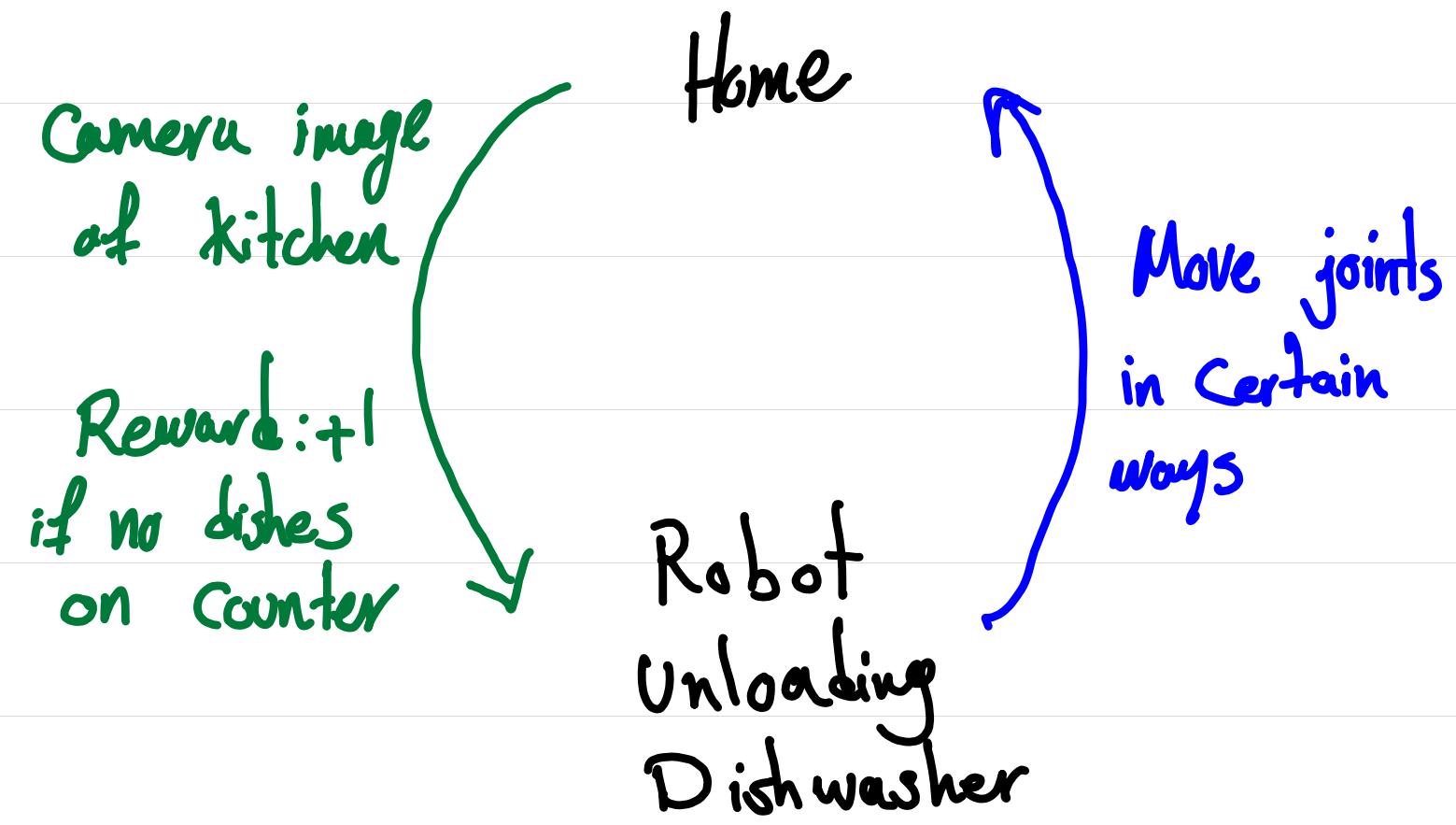
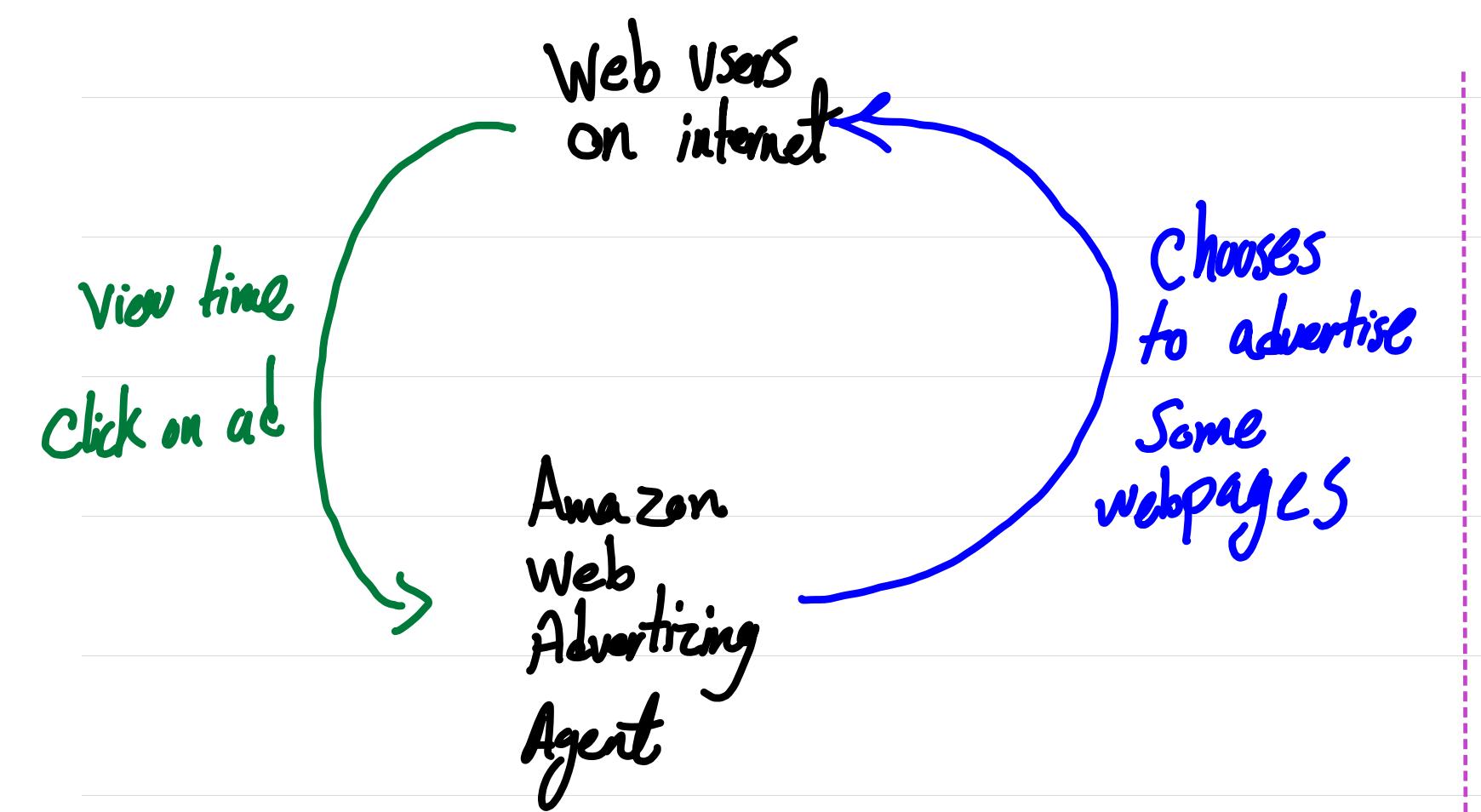
## Sequential Decision Making



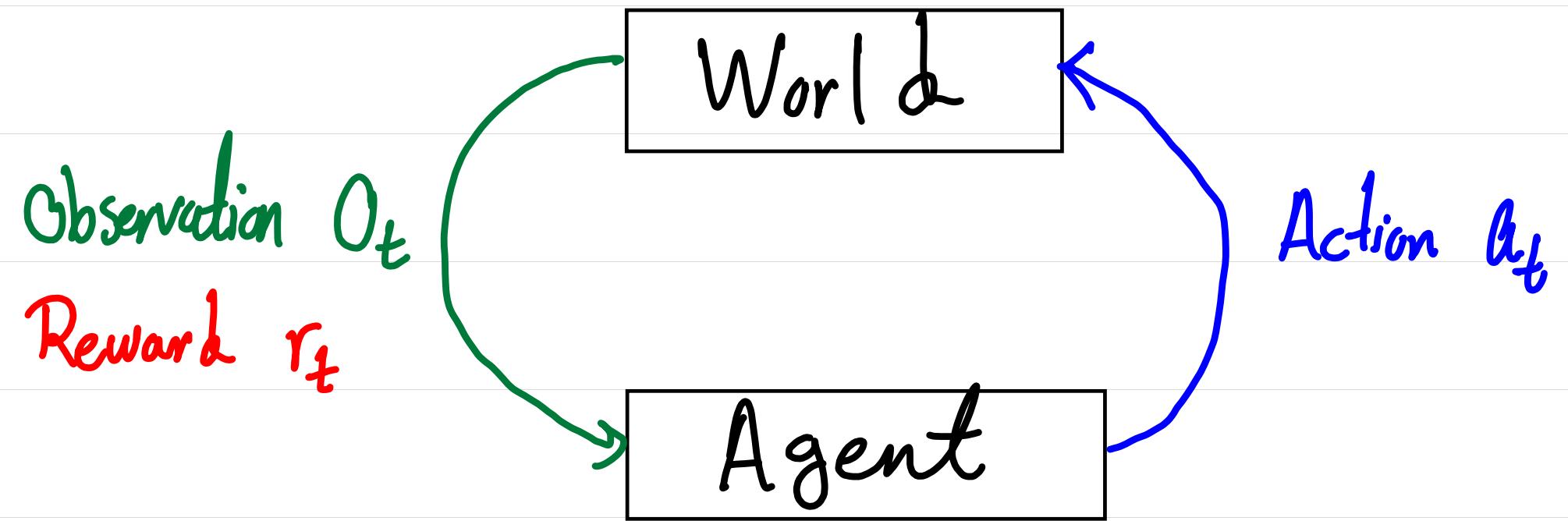
■ Goal: Select actions to maximize total expected future reward

↳ May require balancing immediate & long-term rewards

## Some Examples



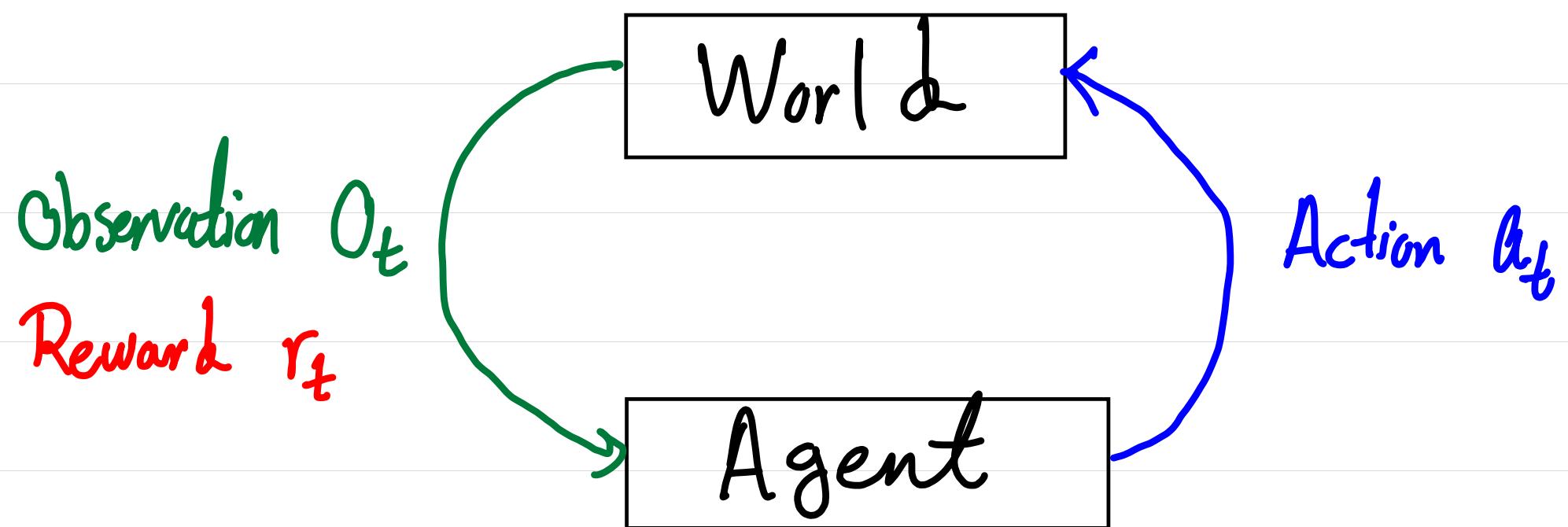
# Sequential Decision Process: Agent & the World (Discrete Time)



- Each time step  $t$ :
  - Agent take an action  $a_t$
  - Given action  $a_t$ , the world updates and



## History: Sequence of Past Observations, Actions, & Rewards



- History:  $h_t = (a_0, o_1, r_1, a_1, o_2, r_2, \dots, a_{t-1}, o_t, r_t)$
- Agent chooses action based on history
- State: information assumed to determine what happens next
  - ↪ State is a Function of history:  $s_t = \Psi(h_t)$

## Observation V.s. History V.s State

- Can anyone tell me what is the difference between observation, history, and state?
- Consider the Atari game. How would you design/choose observation, history, and state?

## Markov Assumption

- Often, to make problems tractable, and because it is not a terrible assumption in reality, we will make **Markov Assumption**.

$$P(S_{t+1} | S_t, a_t) = P(S_{t+1} | S_t, a_t, S_{t-1}, a_{t-1}, S_{t-2}, a_{t-2}, \dots, S_1, a_1)$$

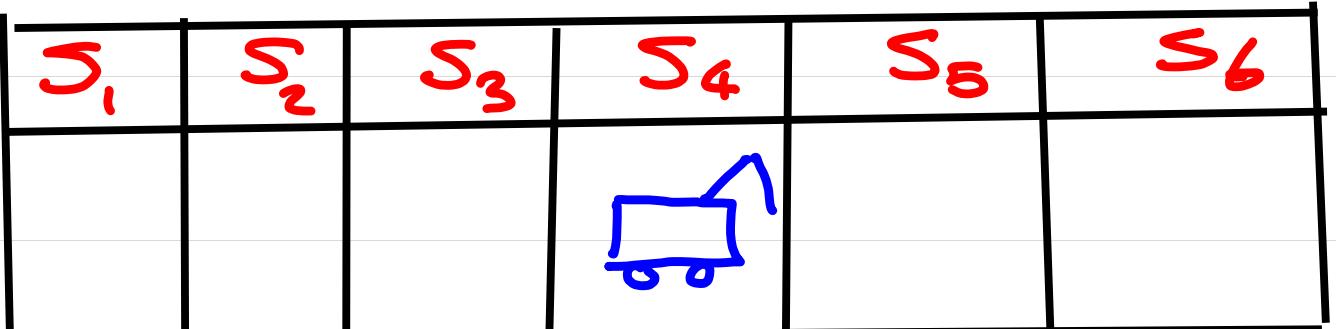
- Future ( $S_{t+1}$ ) is independent of past ( $S_{t-1}, a_{t-1}, \dots, S_1, a_1$ ) given the present ( $S_t, a_t$ )

## Why is Markov Assumption Popular?

- It is simple and often can be satisfied if includes some previous observations as part of state
- There are many problems that we can satisfy Markov assumption by simply setting  $S_t = O_t$ .
- Markov assumption significantly simplifies working with the problem.
- Just like all ML problems we studied so far, there is this trade-off between the expressive power of state representation and how long it takes to train your model.

# MDP Example: Mars Rover

locations on Mars →



- Mars Rover is the robot that you designed to collect samples on Mars.
- Your robot has to choose from going to right or left, to move between locations.
- Mars is strangely small. Only 6 locations on Mars.

State Set:

Action Set:

Rewards:

$s_1$  and  $s_6$  are the terminal locations, you obtain 1 and 5 samples from Soil in locations  $s_1$  and  $s_6$ , respectively, and then the robot will shut down

## Dynamics Model & Reward Model

- Transition / dynamics model: predicts next state
- In other words, dynamics model specifies the distribution of outcomes when the agent makes a decision

- Reward model: predicts immediate reward:
  - ▢ Sometimes just , , or , depending on the model
  - ▢ Note that rewards can be zero, negative, or positive.

## Example: Mars Rover Dynamics Model (Reward Model)

State →

reward →

$s_1$	$s_2$	$s_3$	$s_4$	$s_5$	$s_6$

States & Actions:

Rewards Model:

It's Mars we are talking about. It's different from our Earth. Our robot is not completely predictable there. When the robot decides to go toward a direction, with 0.2 probability, it may go the opposite direction.

## Example: Mars Rover Dynamics Model (Dynamics Model)

State →  
reward →

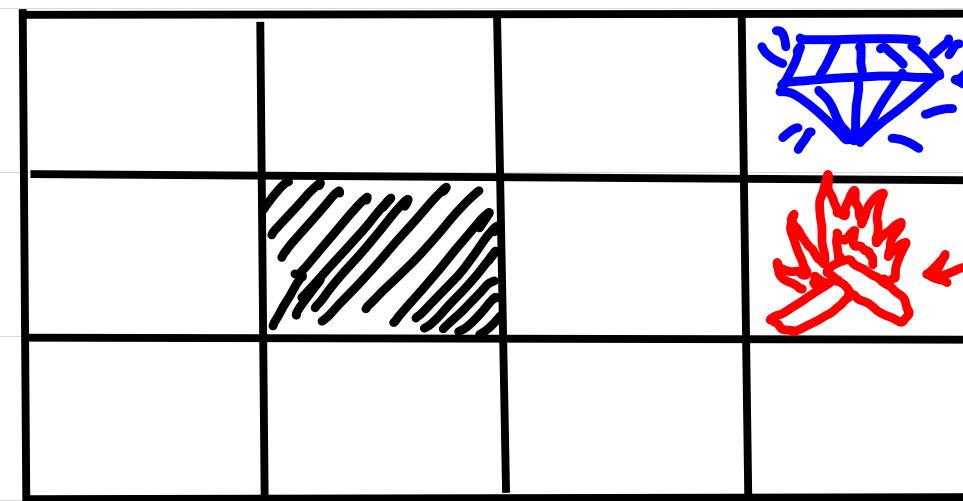
$s_1$	$s_2$	$s_3$	$s_4$	$s_5$	$s_6$

Transition Model:

It's Mars we are talking about. It's different from our Earth. Our robot is not completely predictable there. When the robot decides to go toward a direction, with 0.2 probability, it may go the opposite direction.

## Example: Mars Rover In Grid World

- Consider the following Grid Model of Mars



Mars Rover gets a reward of (+1) if it gets to this location

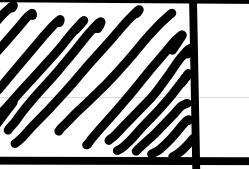
Mars Rover gets a reward of (-1) if it gets to this location

$$R(s, a, s') = -0.03 \text{ for all non-terminal states } s'$$

- Noisy Movements on Mars: actions do not always go as planned
  - With prob. 0.8, the action takes the agent to the desired direction
  - With prob. 0.1, the action takes the agent to a direction perpendicular to the desired direction.
  - If there's a wall in the direction would have been taken to, it stays put.

# Homework: Mars Rover In Grid World

- formulate this MDP.

$S_{2,0}$	$S_{2,1}$	$S_{2,2}$	$S_{2,3}$
$S_{1,0}$		$S_{1,2}$	$2_{1,3}$
$S_{0,0}$	$S_{0,1}$	$S_{0,2}$	$S_{0,3}$

## Example: The undergrad Student's life at UoT

- Erfan wants to thrive in his undergrad
- He has three states: happy, tired, burnt-out
  - burnt-out is the terminal state.
- He has two actions: "Work hard", "Party hard"
  - Working hard gets double reward: reward 2 for working hard



$s$	$a$	$s'$	$T(s, a, s')$	$R(s, a)$
	Work hard		0.5	2
	Work hard		0.5	2
	party hard		1	1
	Work hard		1	2
	Party hard		0.5	1
	party hard		0.5	1
	end		1	0



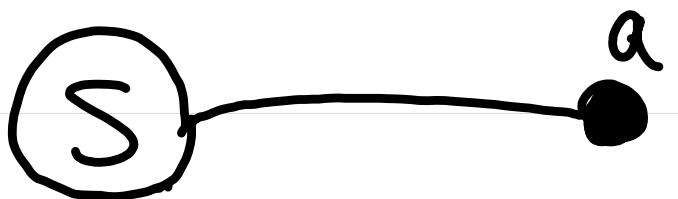
We can summarize  
the dynamics  
and reward model  
with a table

## Transition Graph

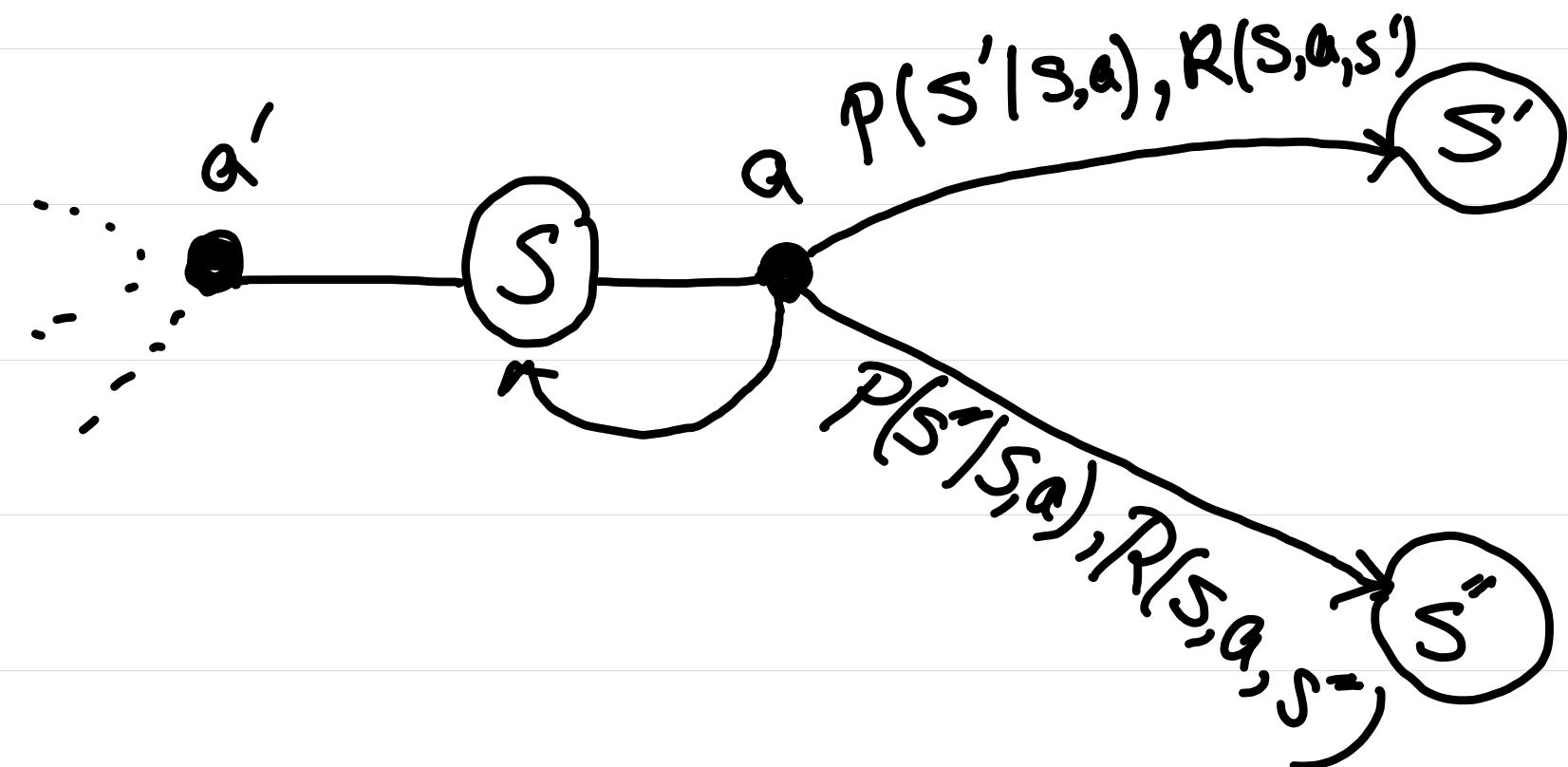
- Transition Graph is a useful way of summarizing the dynamics of a finite MDP
- There are two types of nodes
  - ↳ State nodes: There is a state node for each possible state
    - ▲ Typically a large circle labeled by the name of the state, e.g., 
  - ↳ Action nodes (also called q-state node): There is an action node for each state-action pair.
    - ▲ Typically a small solid circle labeled by the name of action or (state, action) pair and connected by a line to the state node, e.g.,  or 

## Transition Graph

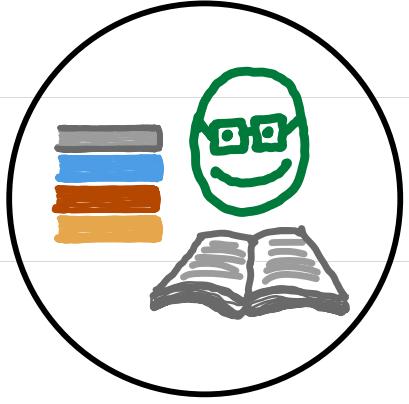
- Starting in state "s" and taking action "a" move you along the line from state nodes "s" to action node  $(s, a)$ .



- Then the environment responds with a transition to the next state's node.

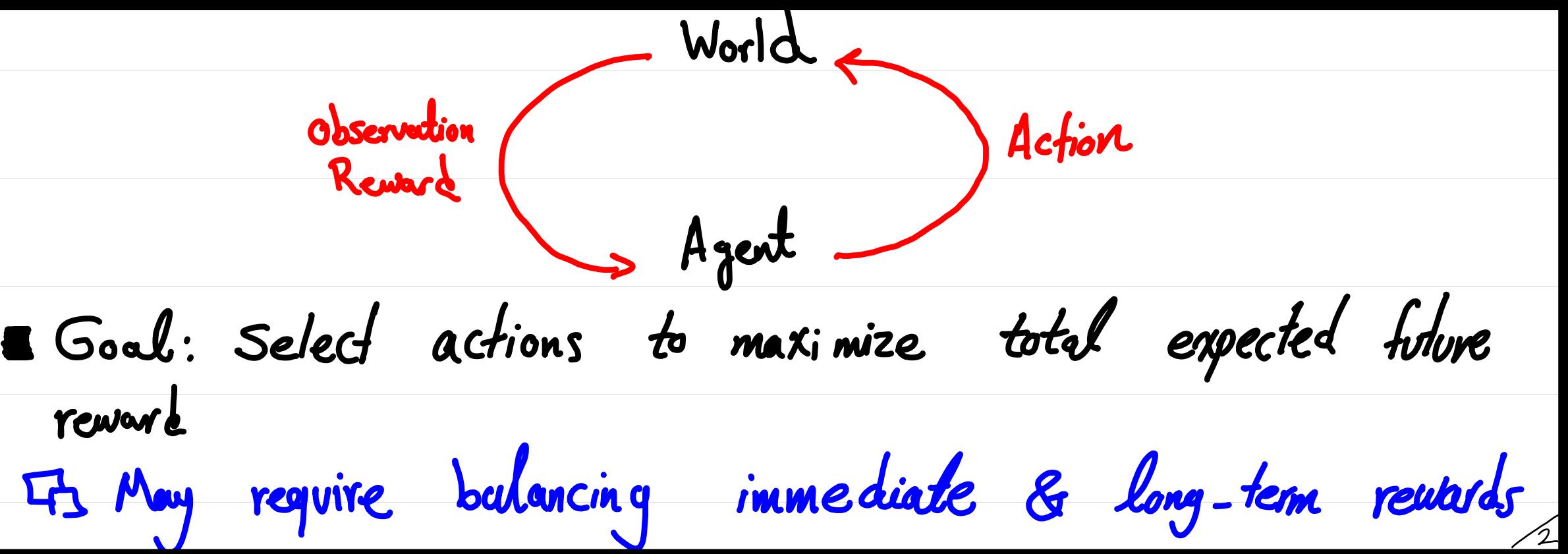


Activity: Draw the transition Graph for Erfan the Undergrad



# Returns & Episodes

■ Recall:



■ Goal: Select actions to maximize total expected future reward

↳ May require balancing immediate & long-term rewards

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

↳ Return is a function of the reward sequence

## 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  $\infty$ .

## 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

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

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

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

■ As  $\gamma$  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  $\pi$ , 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*  $\pi^*$ , 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  $a_e$ .

## 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 :

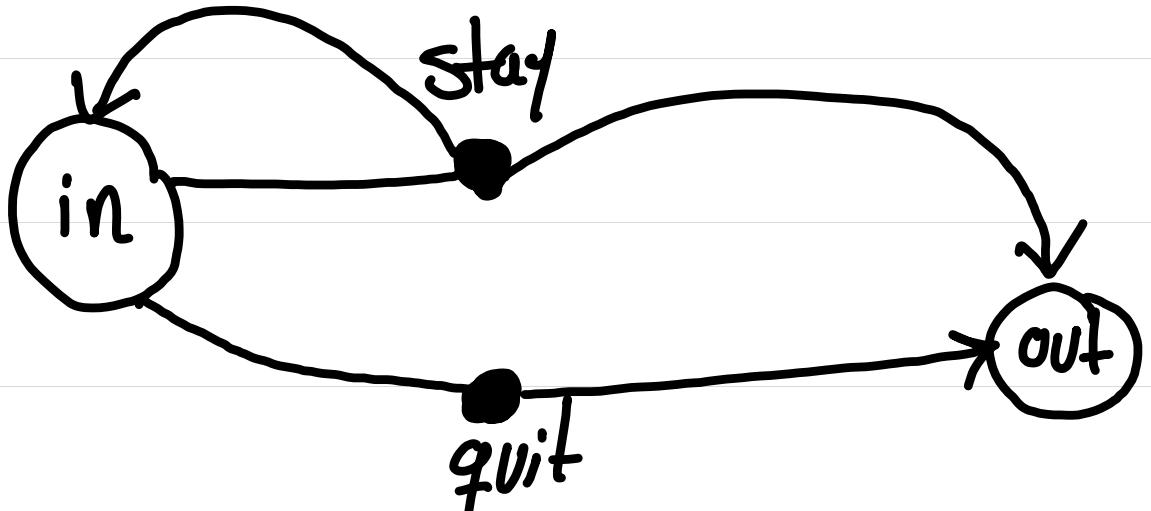
Erfan the gambler  
(See the next page for the description)

## 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

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

8

:

:

.

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

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

■ Starting at state " $s$ " and following policy " $\pi$ " 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  $\pi$ . 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

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

- 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  $R(s, a)$

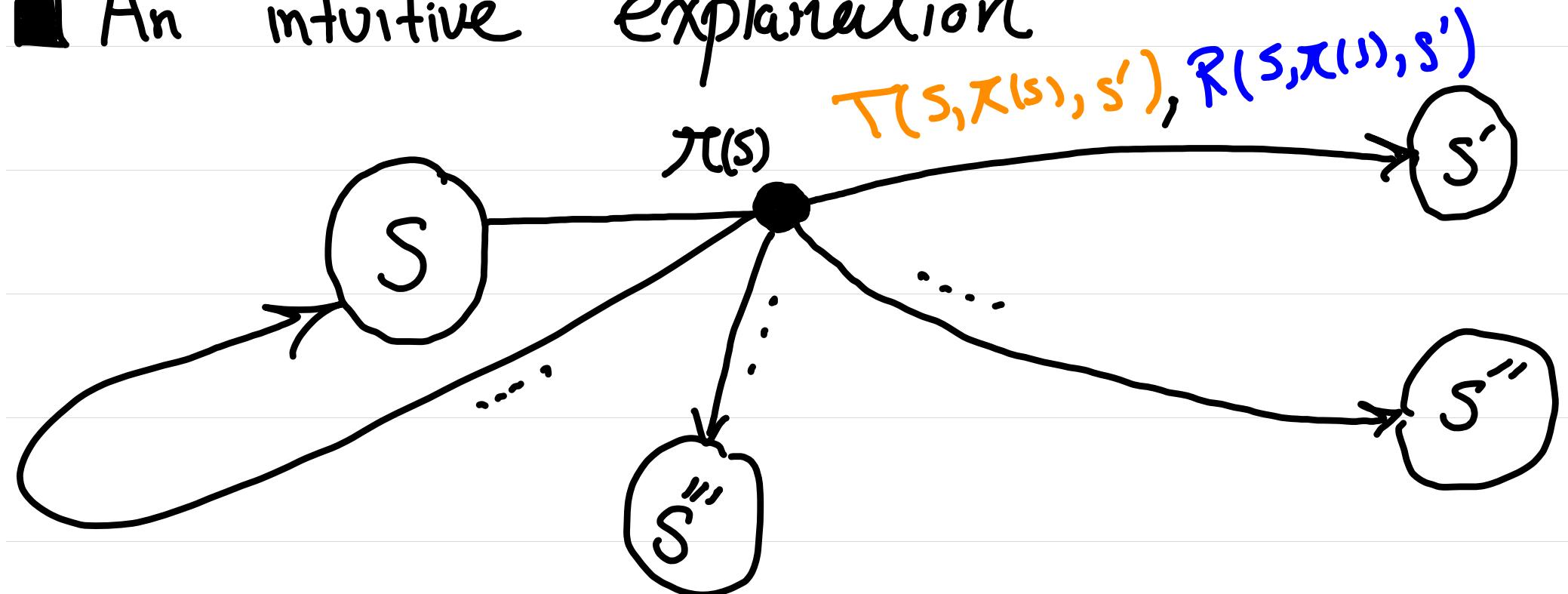
$$V_\pi(s) = R(s, \pi(s)) + \gamma \sum_{s'} T(s, \pi(s), s') V_\pi(s')$$

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

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

## 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  $\gamma$ .

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

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

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

## Value and Q-value of a Policy

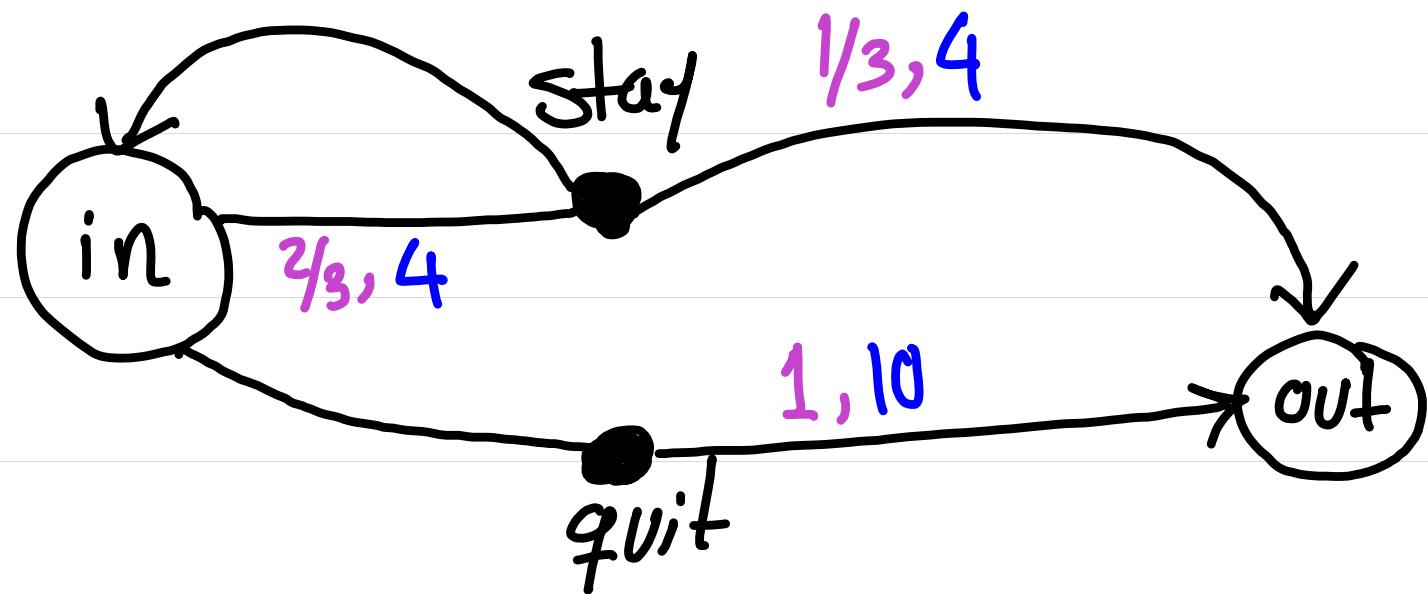
- We use  $V_\pi(s)$  to denote the *expected utility* received by following policy " $\pi$ " from state "S".
- We use  $Q_\pi(s, a)$  to denote the expected utility of taking action "a" from state "s", and then following policy " $\pi$ ".

$V_\pi(s) \leftarrow$  value of a policy

$Q_\pi(s, a) \leftarrow$  Q-value of a policy

## Example: Erfan the gambler

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



- Let  $\pi_i$  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  $\pi$ .

■ Observe that

$$V_{\pi_i}("in") = \sum_{s'} T("in", \pi_i(in), s') \left[ R(in, \pi_i(in), s') + \gamma V_{\pi_i}(s') \right]$$

=

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

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

$$V_{\pi_2}(\text{in}) = \sum_{s'} T(\text{in}, \pi_2(\text{in}), s') [R(\text{in}, \pi_2(\text{in}), s') + \gamma V_{\pi_2}(s')]$$

=

■ Which Policy do you think 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,  $\underline{V} = (I - \gamma T)^{-1} ((T \odot R) \underline{1})$



