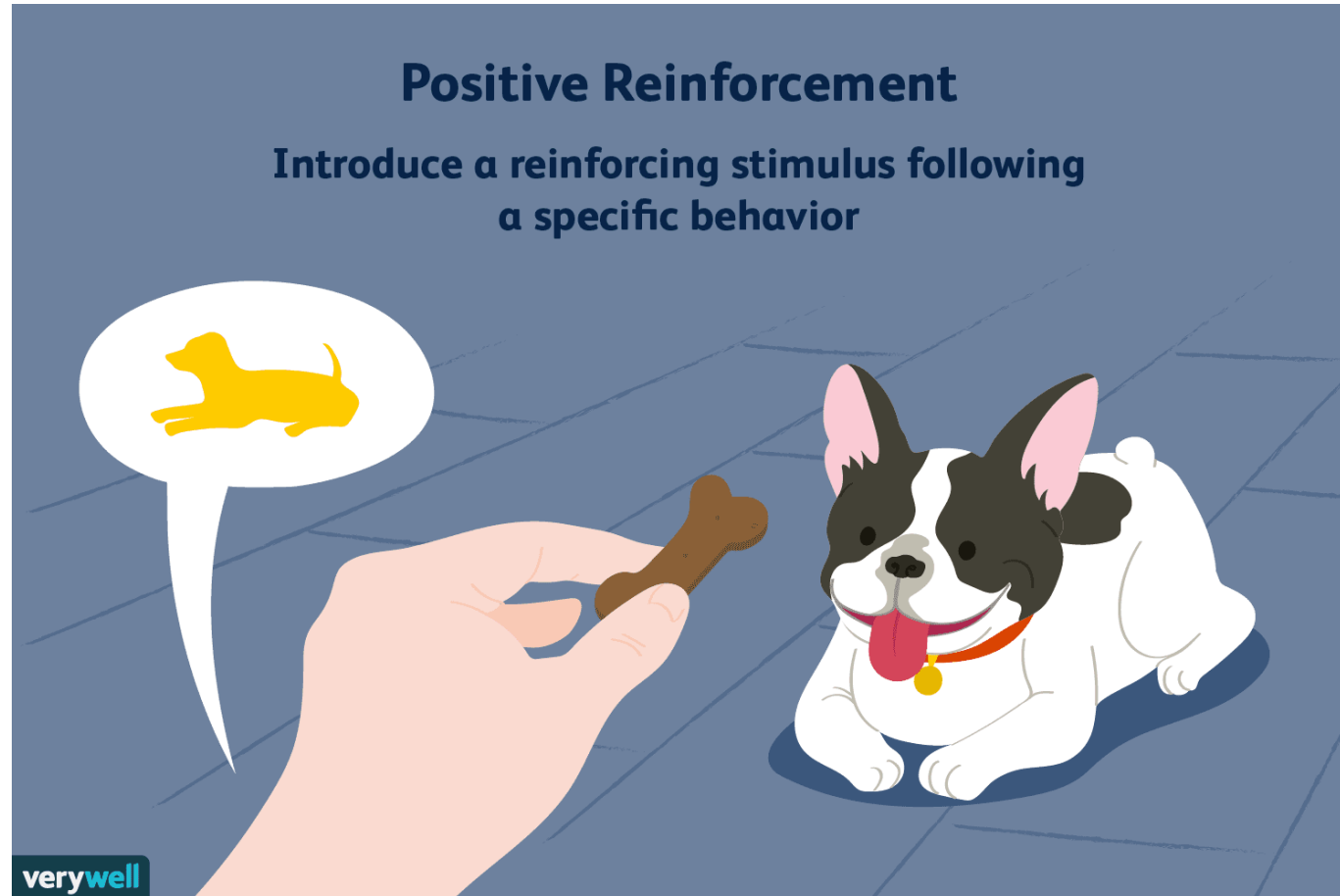


Markov Property & Markov Decision Process Reinforcement Learning Review

Based on Prof. Oh's Reinforcement Learning Lectures

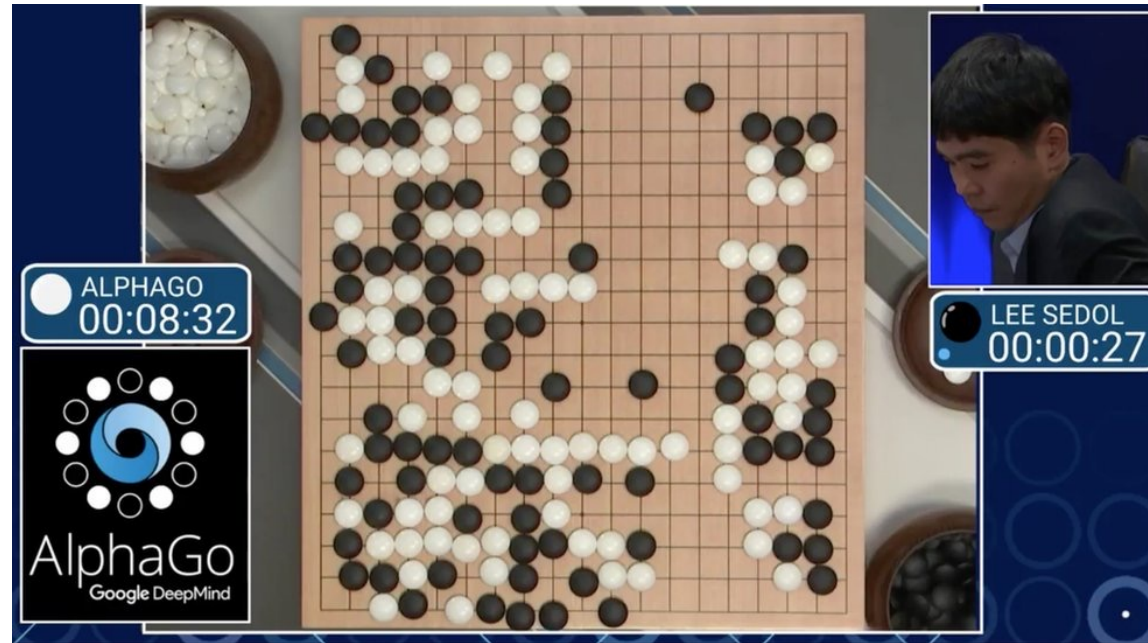
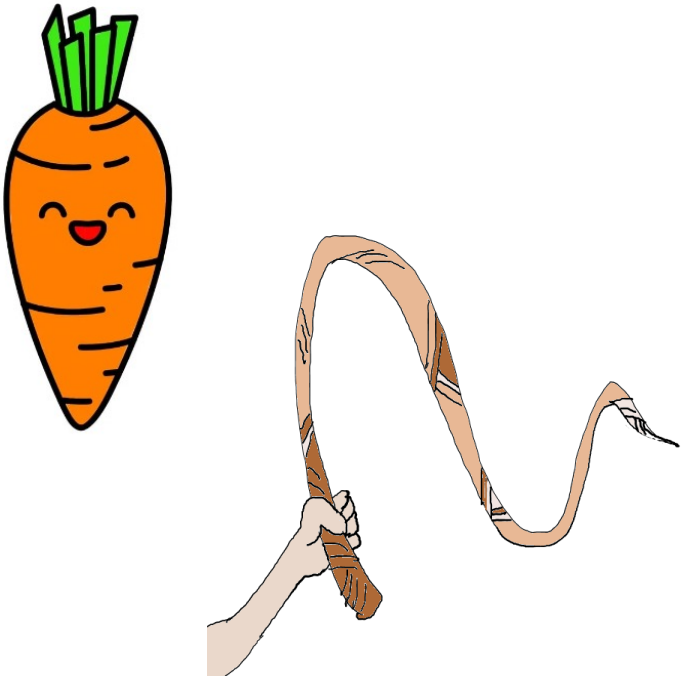
Suinne Lee

Reinforcement Learning

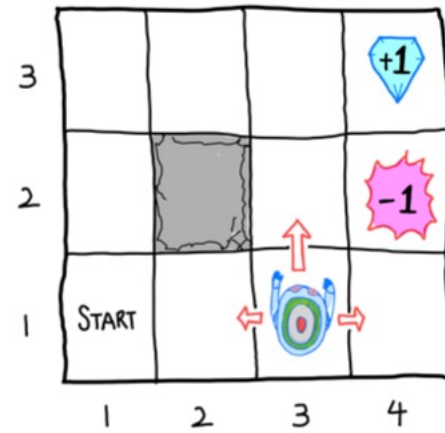


Train something to be good at something through trial and error

More suited for something where the agent actively engages with the environment



Example: Grid World

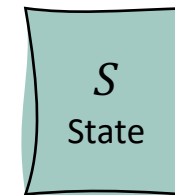


The Very Basics

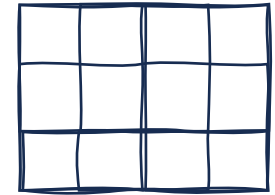
- Agent
 - The player? Robot? Something that takes action!
- State:
 - Environment? Setting? Surrounding?
 - Capital S: notation for the set of all possible states, or *state space*
- Action
 - What is done
- Reward
 - Something as a measure of how good the state/action was.

Goal: maximize reward!

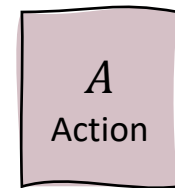
In the Grid world Example:



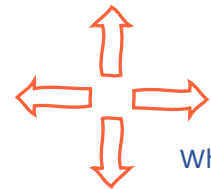
=



The location
(one of these squares!)



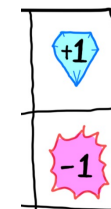
=



Which direction to move

R
Reward

=



These (good squares to land
/ bad squares to land)
+ small negative reward for
path length

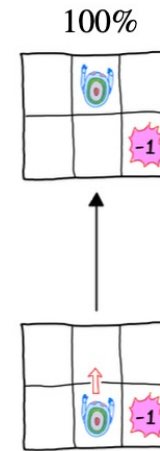
- State \rightarrow Choose Action. (Policy)

- Action \rightarrow Next State (state transition, given action)

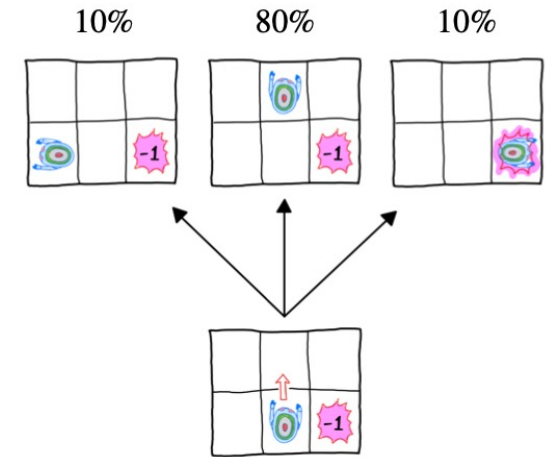
Each of these steps may be noisy (stochastic)

- State \rightarrow Next State (state transition)

Deterministic grid world



Stochastic grid world



So we also establish the following concept:

- Policy

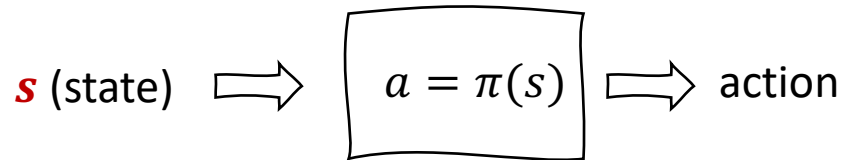
- Deterministic: $\pi(s) = a$

- Function of **state**. Gives a determined function value: which is the action.

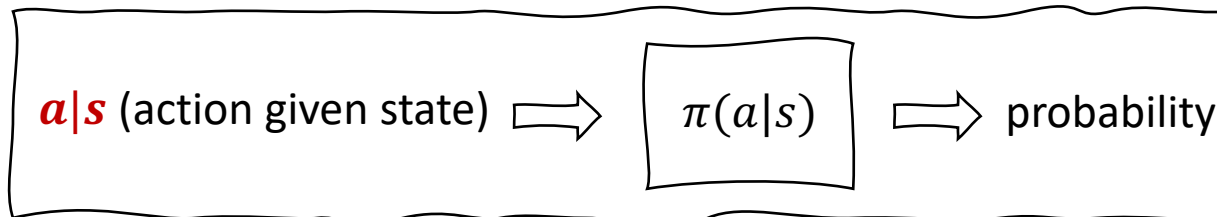
- Stochastic $\pi(a|s)$

- **Probability distribution** of **state**. Which means the notation $\pi(a|s)$ itself means a completely different thing than $\pi(s) = a$: $\pi(a|s)$ in itself means the **probability mass or density**, as a function of the **variable $a|s$** (action given state)

Deterministic:



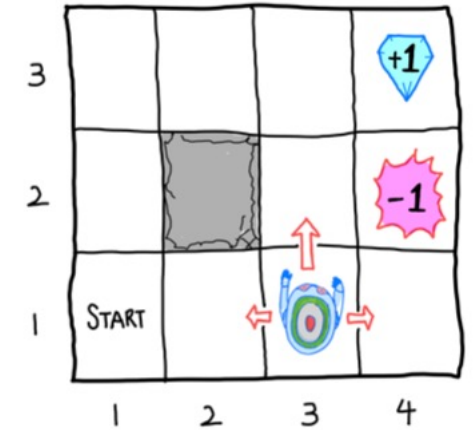
Stochastic:



As a whole, represents action taken given a state in the form of a probability distribution.

Back to the Grid World Example

Episode $(1, 1) \xrightarrow{\text{north}} (1, 2) \xrightarrow[\text{wall}]{\text{east}} (1, 2) \xrightarrow{\text{north}} (1, 3) \xrightarrow[\text{east 10\%}]{\text{south}} (2, 3) \xrightarrow{\text{east}} (3, 3) \xrightarrow{\text{east}} (4, 3)$
with the total reward $5c+1$



Episode: one full game. 한 판.

(a bit of a vague term, but we can also say that if the state, state transitions (and thus rewards) are identical for two episodes, they are the same episode.)

Taking this into account, depending on whether the state transition is deterministic or stochastic, a single policy may bring about one or multiple episodes.

Markov Property

If we were to consider the state transition probability in a general model, the future state would depend on all past states as well as the present state. But the Markov Property restricts this to a specific case where it doesn't.

- $P(S_{t+1}|S_0, S_1, \dots S_t) = P(S_{t+1}|S_t)$
- The future state does not depend on the past states (independent of them), only on the present state.
- In information theory:

- Random variables X, Y, Z are said to form a **Markov chain** in that order (denoted by $X \rightarrow Y \rightarrow Z$) if the conditional distribution of Z depends only on Y and is conditionally independent of X
- Specifically, $X \rightarrow Y \rightarrow Z$ if

$$p(x, y, z) = p(x)p(y|x)p(z|y)$$

Some simple consequences

- $X \rightarrow Y \rightarrow Z$ if and only if X and Z are conditionally independent given Y

$$p(x, z|y) = \frac{p(x, y, z)}{p(y)} = \frac{p(x, y)p(z|y)}{p(y)} = p(x|y)p(z|y)$$

- $X \rightarrow Y \rightarrow Z$ implies $Z \rightarrow Y \rightarrow X$
- If $Z = f(Y)$, then $X \rightarrow Y \rightarrow Z$

Markov Process (Markov Chain)

- Stochastic process:

- def: collection of random variables indexed by time (can be generalized to a continuous case)

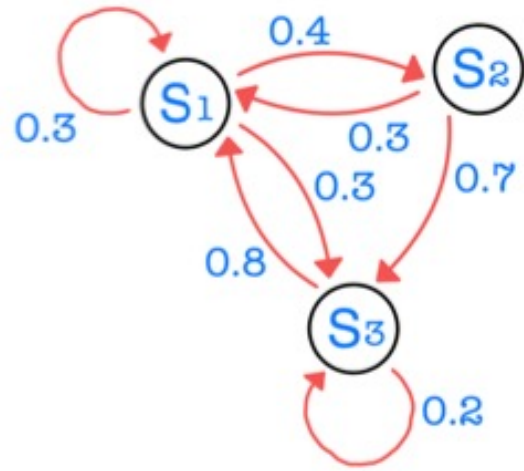
$$S_0, S_1, \dots, S_{t-1}, S_t, S_{t+1}, \dots$$

- The Markov Process is a process where the Markov Property holds.
- In RL, we assume the processes to be Markov.

Markov Process: Mathematical Definition

- Defined as the tuple (S, P)
 - S : States. (Variables, not specific values!)
 - P : Transition matrix
$$P_{ij} = P_{s_i s_j} = p(s_j | s_i) = P(S_{t+1} = s_j | S_t = s_i)$$
 - Naturally, the sum of elements in each row (outgoing from state s_i) is 1.

Markov Process



	S_1	S_2	S_3
S_1	0.3	0.4	0.3
S_2	0.3	0.0	0.7
S_3	0.8	0.0	0.2

Markov Decision Process (MDP)

- Elementary probabilistic model for implementing RL

Markov Decision Process (MDP)

- S: State space
- A: Action space

State and action spaces need not be finite.

- P: transition probability, from s to s' given a

$$P_{ss'}^a = p(s' | s, a) = P(S_{t+1} = s' | S_t = s, A_t = a)$$

- R: Reward function

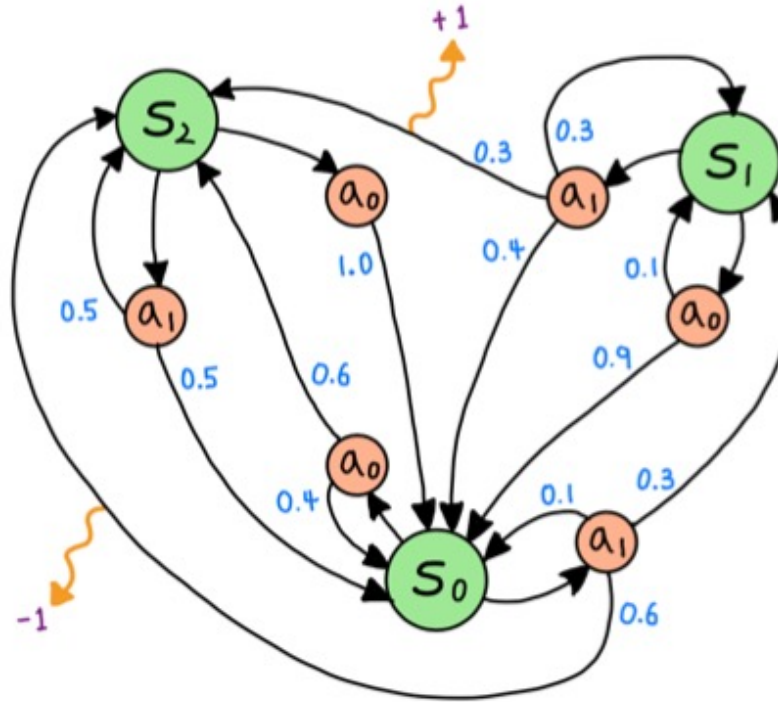
$R_{ss'}^a$ (possibly 3 different types: R_s , R_s^a , $R_{ss'}^a$)

The reward may depend on

- the current state
- the current state and action taken
- the current state, action, next state achieved

- Discount factor $\gamma \in [0, 1]$: to be explained (has to do with prioritizing immediate rewards)

Markov Decision Process (MDP)



3 states: S_0, S_1, S_2

2 actions: a_0, a_1

2 rewards: $+1, -1$

$$P_{S_0 S_1}^{a_1} = ?$$

$$R_{S_1 S_2}^{a_1} = ?$$

Example → next talk!

Glossary

- Agent
- State
- Action
- Reward
- Policy
- Episode
- Markov Property
- Markov Chain / Markov Process