# REINFORCEMENT LEARNING

→ Richard Sutton. & Barto. :- http://incompleteideas.net/book/the-book-2nd.html
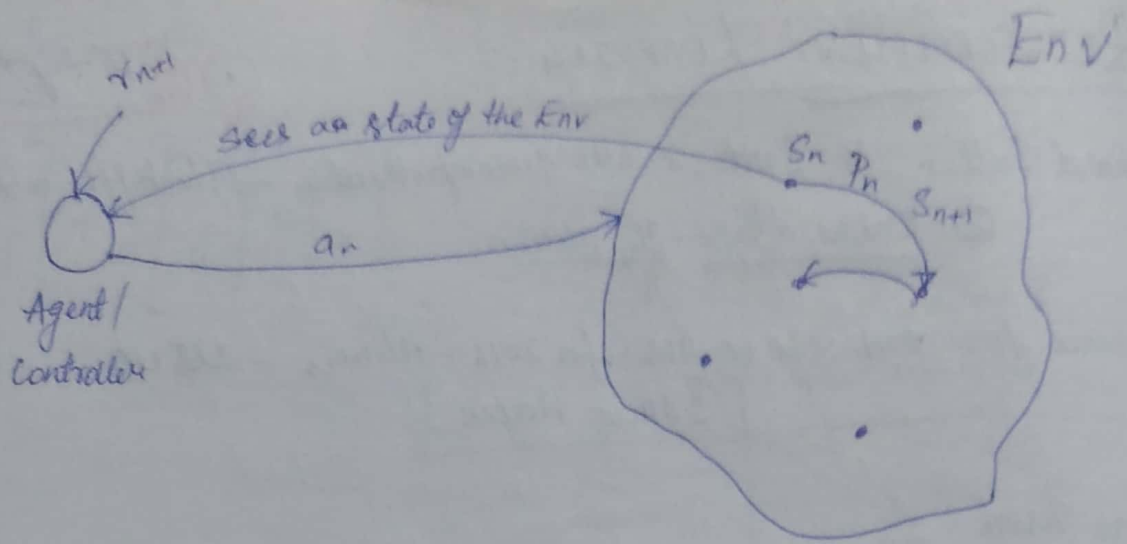    Online draft - II edition

→ Reference for prob : ece 313 - fa 2016 - illinois - UIUC
                [Bruce Hajek]

Coure Eval :
    → Project evaluation
    → Mid term
    → Final Exam
    → Assignments.



Learning

Supervized Learning
(teacher-Student)

Unsupovized learning
(trail - Error)

R.L

Env

sees an state of the Env

Agent / Controller

$S_n$ — State at time 'n'

$a_n$ — action / control at time 'n'

$r_n$ — reward at time 'n'
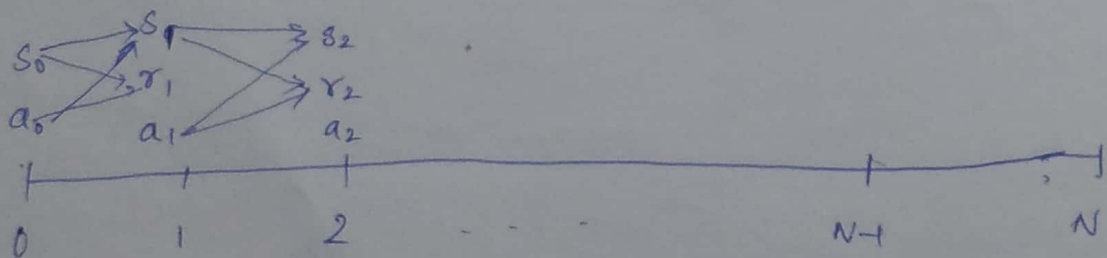
$P_n$ — Prob that env moves to state $S_{n+1}$
(given that current state = $S_n$ & agent picks control / action $a_n$)

## Important things to note:

1. Dynamic decision making

2. Uncertain environment.

## Goal of the agent:

Learn a sequence of actions or controls given the date of the environment in order to <u>Max</u> long term reward.



X. Agent is oblivious to the model of the system

$$\max \; E\left[ \sum_{n=0}^{N} r_n \right] \quad —\text{①} \qquad \longleftarrow \text{ if the time horizon is finite.}$$

When $N = \infty$,

eqn① doesn't fit since $\overset{\infty}{\underset{0}{\sum}}$ is always towards maxima.

Hence, discounted cost problem

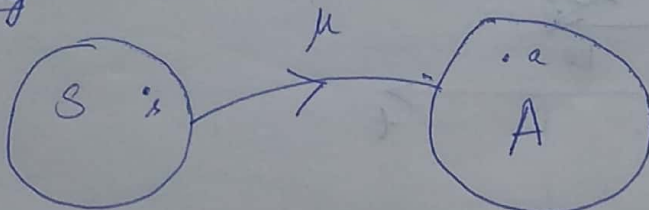$$\max \; E\left[ \sum_{n=0}^{\infty} \gamma^n r_n \right] \qquad \qquad E - \text{Expected Value.}$$

Where discount factor $\gamma \in (0,1)$

## Exploration vs Exploitation

Need to judiciously combine ~~Expl~~ explor$^n$ & exploit$^n$


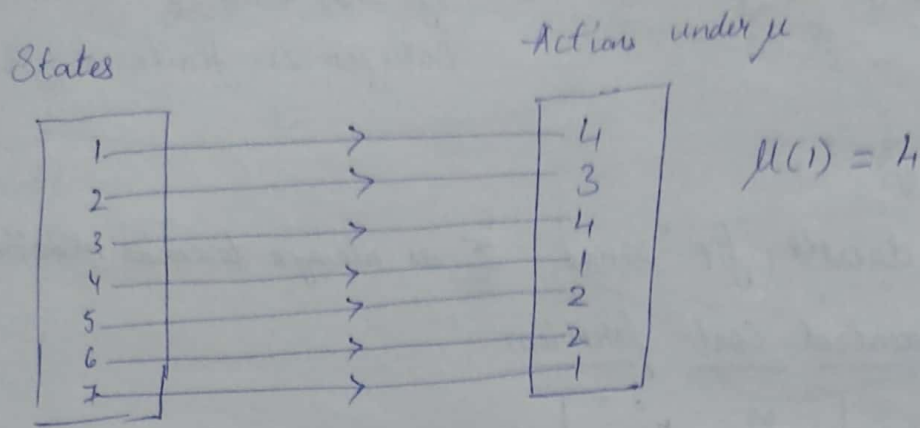
1. **Policy:** A function from states to action



$$\mu(s) = a$$

$S$ : set of all states

$A$ : set of all actions.

Kavitha

Eg: States = 7, action = 4, considering a deterministic policy μ

States                   Action under μ



$\mu(1) = 4$

**probabilistic actions:**
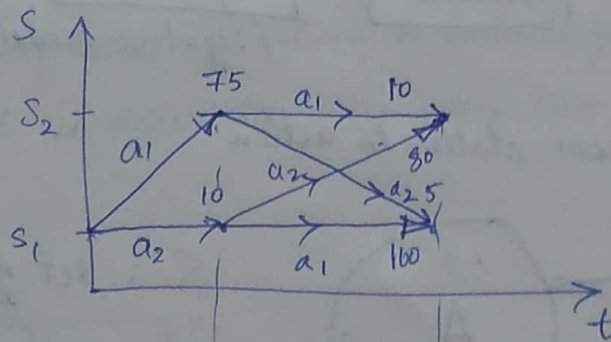
$\phi(s, 1) = \emptyset$
$\phi(s, 2) =$

$\phi = \phi(\phi(s,a), s \in S, a \in A)$

## 2. Reward function

A map that inputs state action tuples & gives a numerical value.

→ Reward in a state is indicator of "immediate desirability" of the state.



$a_1 @ 0$           $a_1 @ 0$        ∴ long term vs

is prefered        is prefered under μ     dan short term desirability

# 3. Value function

A map that inputs states or state-action tuples and output a numerical quantity.

- Value function tells us the "_long-term_" desirability

- Decision making will involve _max_ value $f^n$

- A state may get a _low reward_ yet have _high value_


# 4. Model of the environment

→ model emulates the environment

→ Random - state transitions (transition prob $s_i \to s_j$)

$$p(S_1, a_1, S_2)$$

(the randomness is because of $s_i \to s_j$ is not fixed)

# Temporal difference learning methods

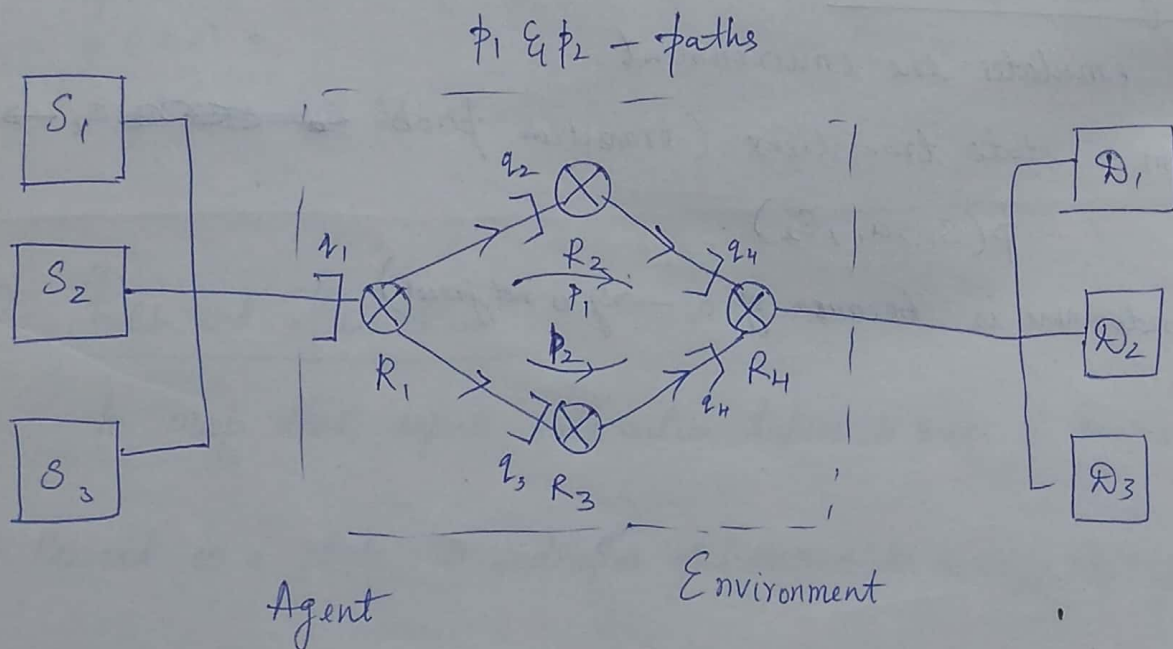For a given policy (fixed for the entire time duration), estimate the value function.

$$V(s) \leftarrow v(s) + \alpha \left[ v(s') - v(s) \right]$$

(new estimate) of value of $s$

$s \longrightarrow$ current state

$s' \longrightarrow$ next step.

$\alpha \longrightarrow$ a small number (increment)

$p_1 \& p_2$ — paths



Agent            Environment

$$S = (q_1, q_2, q_3, q_4)$$

$q_i = \underset{=}{no}$ of packets at router $i$, $i = 1, 2, 3, 4$.

$$a = (p_1, p_2) \in \{ (1, 0) \text{ or } (0, 1) \}$$

         either $p_1$ or $p_2$ can be take not both at same time.

Objective : minimize delay

$$p_1 \longrightarrow -(q_1 + q_2 + q_4)$$

$$p_2 \longrightarrow -(q_1 + q_3 + q_4)$$

( negative state
have least of no of
packets in the queue
waiting )

$$\max_{\{p_1, p_2\}} \mathbb{E}\left[\sum_{n=0}^{\infty} r^n g_n\right]$$

---

## Tic - Tac - Toe

**Goal:** learn optimal policy for player given that opponent plays same policy.

**State description :-**

$$s = \{S_1, S_2, S_3 \cdots S_9, S_{10}\}$$

| $S_1$ | $S_2$ | $S_3$ |
|-------|-------|-------|
| $S_4$ | $S_5$ | $S_6$ |
| $S_7$ | $S_8$ | $S_9$ |

$$S_1, S_2 \cdots S_9 \in \{e, x, 0\}$$

$\longleftarrow$ who played first

$$S_{10} \in \{x, 0\}$$

e — empty
x — player 1
0 — player 2

$$\text{Reward} = \begin{cases} 1 & \text{if player has won} \\ 0 & \text{otherwise} \end{cases}$$

The game can be represented as tree structure,

Tree structure (suppose opponent is first player) & $S_{10} = 0$



exploration

**Q₁:** Suppose both player & opponent use same R.L algo to learn their moves.

**Case 1:** Reward $\begin{cases} 1 & \text{if Player/opp wins} \\ 0 & \text{if drawn} \\ -1 & \text{if player/opp looses} \end{cases}$  Say $S_{10}$ - random

Zero sum game,

In case 1, since the rewards are symmetric & thus they learn same policy. & learning will converge.

**Case 2:** Reward $\begin{cases} 1 & \text{if player to wins} \\ 0 & \text{other wise} \end{cases}$  Say $S_{10}$ - random

in this case, learning might not converge, as the rewards are is not zero sum. Since both the players try to outsmart other unlike typical RL setup one player against Env

so for two player RL games, symmetric rewards ensure convergence.

$Q_2$: Should symmetrically eqt positions have same values

A: Yes



$$V(S_1) = V(S_2)$$

$Q_3$: Greedy play: Suppose RL player is greedy, i.e., always in exploit mode, What are the problems?

Ans:
→ smarter opponents will affect in future might affect player

→ policy learnt might be the best as the different rewards are not explored.

Kavitha