

## Action-Values

- \* The value is the expected reward

$$Q_t(a) \doteq E[R_t | A_t = a] \quad \forall a \in \{1 \dots k\}$$

$$= \sum_r p(r|a) r$$

- \* The goal is to maximize the expected reward

argmax  $Q_t(a)$

- \* Calculating  $Q_t(a)$

$$\frac{1}{-11} \quad | \quad \frac{1}{9} \quad | \quad \frac{1}{9} \rightarrow Q_t(a) = -0.5 \times -11 + 0.5 \times 9$$

↓      ↓      ↓      ↓  
 prob of reward when failure    prob of reward when success  
 fail      success      success      success

$\rightarrow Q_t(a)$  is not known

Sample-Average Method

$\hat{Q}_t(a) = \frac{\text{sum of rewards when a taken prior to } t}{\text{number of times a taken prior to } t}$

$$= \frac{\sum_{i=1}^{t-1} R_i}{t-1}$$

## Action Selection

greedy action

$$a_g = \arg\max Q(a)$$

Summary :-

- Sample average method can be used to estimate action values.
- The greedy action is the action with the highest value estimate.
- \* Incremental update rule

$$Q_{n+1} = \frac{1}{n} \sum_{i=1}^n R_i$$

pull current reward out of the sum

$$= \frac{1}{n} \left( R_n + \sum_{j=1}^{n-1} R_i \right)$$

$$= \frac{1}{n} \left( R_n + (n-1) \frac{1}{n-1} \sum_{j=1}^{n-1} R_i \right)$$

$$= \frac{1}{n} (R_n + (n-1) Q_n)$$

Recall (current value estimate)

$$\hat{Q}_n = \frac{1}{n-1} \sum_{i=1}^{n-1} R_i$$

$$= \frac{1}{n} (R_n + (n\hat{Q}_n - \hat{Q}_n))$$

$$\hat{Q}_{n+1} = \hat{Q}_n + \frac{1}{n} (R_n - \hat{Q}_n)$$

NewEstimate  $\leftarrow$  OldEstimate +  
StepSize [Target - OldEstimate]

$$\hat{Q}_{n+1} = \hat{Q}_n + \alpha_n (R_n - \hat{Q}_n)$$

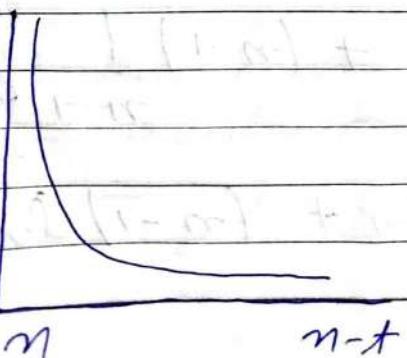
$$\alpha_n \in [0, 1]$$

$$\alpha_n = \frac{1}{n}$$

Non-stationary bandit problem :-

The distribution of rewards changes with time

Reward  
Weight



\* Decaying past rewards

$$Q_{n+1} = Q_n + \alpha_n (R_n - Q_n)$$

$$= \alpha_n R_n + Q_n - \alpha_n Q_n$$

$$= \alpha_n R_n + (1-\alpha) Q_n$$

$$= \alpha R_n + (1-\alpha)[\alpha R_{n-1} + (1-\alpha) Q_{n-1}]$$

$$= \alpha R_n + (1-\alpha)\alpha R_{n-1} + (1-\alpha)^2 Q_{n-1}$$

$$= \alpha R_n + (1-\alpha)\alpha R_{n-1} + (1-\alpha)^2 \alpha R_{n-2} + \dots + (1-\alpha)^{n-1} \alpha R_1 + (1-\alpha)^n Q_1$$

$$= (1-\alpha)^n Q_1 + \sum_{i=1}^n \alpha (1-\alpha)^{n-i} R_i$$

$Q_1$  > initial action value

\* The first term tells us that the contribution of  $Q_1$  decreases exponentially with time.

\* The second term tells us the rewards further back in time contribute exponentially less to the sum.

The most recent rewards contribute most to our current estimate

\* Exploration & Exploitation trade-off

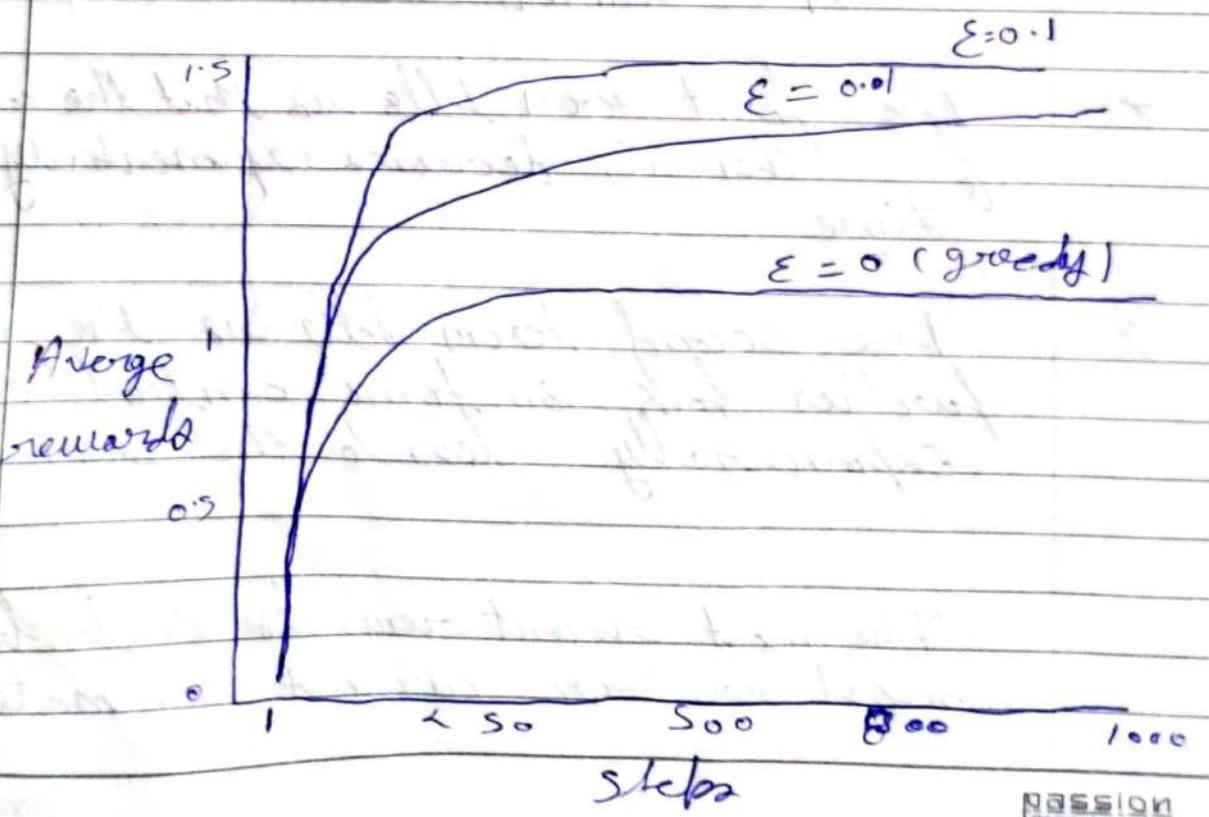
Exploration - improve knowledge for long term benefit

Exploitation - exploit knowledge for short term benefit

→ Epsilon-Greedy Action Selection

epsilon - probability of choosing to explore

$A_t \leftarrow \begin{cases} \text{argmax}_a Q_t(a) & \text{with prob } \epsilon \\ a \sim \text{Uniform}(a_1, \dots, a_K) & \text{with prob } 1 - \epsilon \end{cases}$



### Clinical Trials :-

A reward of 1 if the treatment succeeds  
otherwise 0

Initial Values  $\rightarrow$

$$\begin{aligned} \theta_1(P) &= 2.0 \\ \theta_1(Y) &= 2.0 \\ \theta_1(B) &= 2.0 \end{aligned}$$

$$\theta_{n+1} = \theta_n + \alpha (R_n - \theta_n)$$

Let  $\alpha = 0.5$

doctor choose P

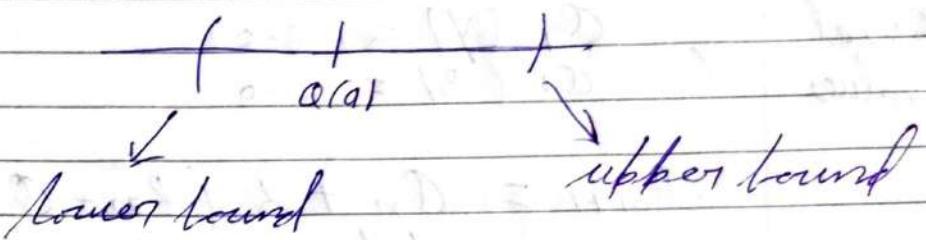
$$\begin{aligned} \theta_2 &= 2 + 0.5(1 - 2) \\ &= 2 - 0.5 \\ \theta_2 &= 1.5 \end{aligned}$$

Limitations of optimistic initial values

- ① optimistic initial values only derive early exploration
- ② They are not <sup>well</sup> suited for non-stationary problem.
- ③ We may not know the what the optimistic initial value should be.

## \* Upper-Confidence Bound (UCB) Action selection

Uncertainty in estimates



- The region in blue is the CI which represents over uncertainty.
- If this region is large, we are uncertain that the value of action A is near or estimated value.
- If region is very small, we are very certain that the value of action A is near estimated value.

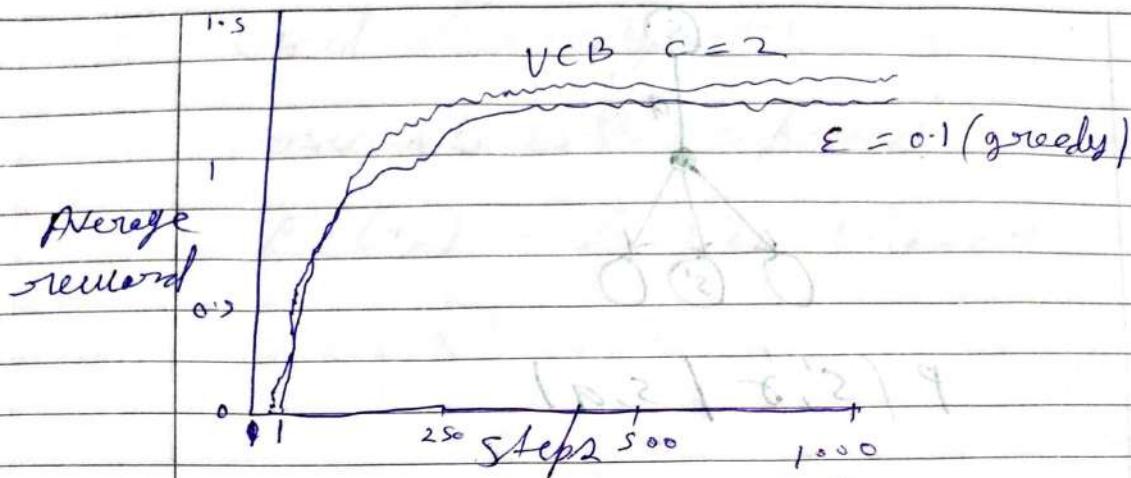
$$A_t = \arg\max [Q_t(a) + C \sqrt{\frac{\ln t}{N_t(a)}}]$$

↑ exploit                      ↑ explore

$C \rightarrow$  control the amount of exploration

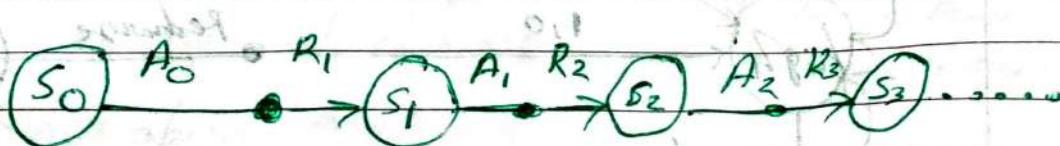
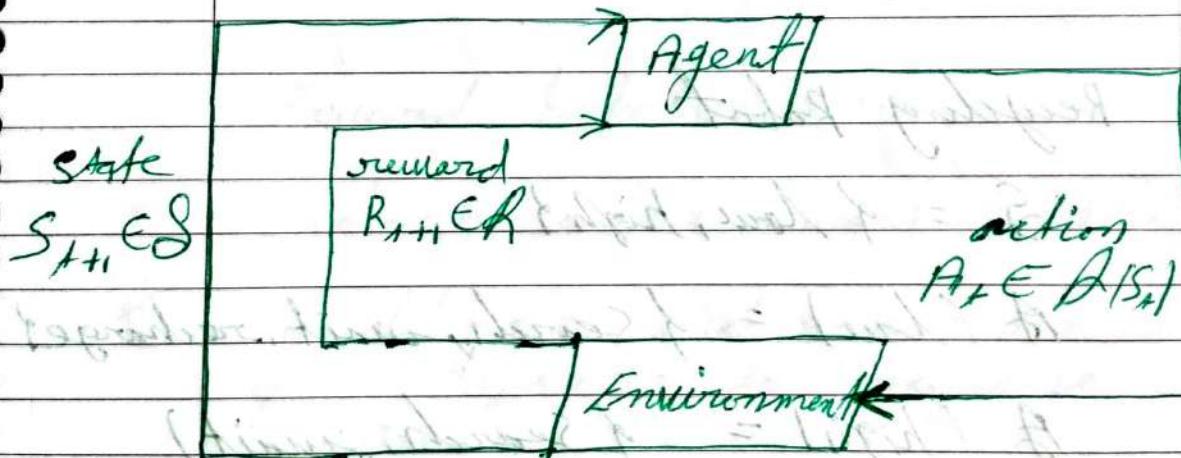
$$C \sqrt{\frac{1.2 \text{ timesteps}}{\text{times action } a \text{ taken}}}$$

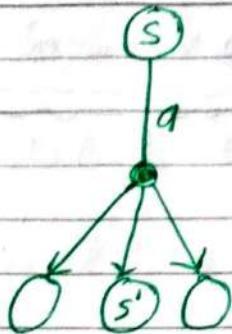
Date - / - / -



course - 1 week - 2

## Markov Decision processes (MDP)





$$P(S', \sigma | s, a)$$

$$\sum_{S' \in S} \sum_{\sigma \in \Sigma} P(S', \sigma | s, a) = 1$$

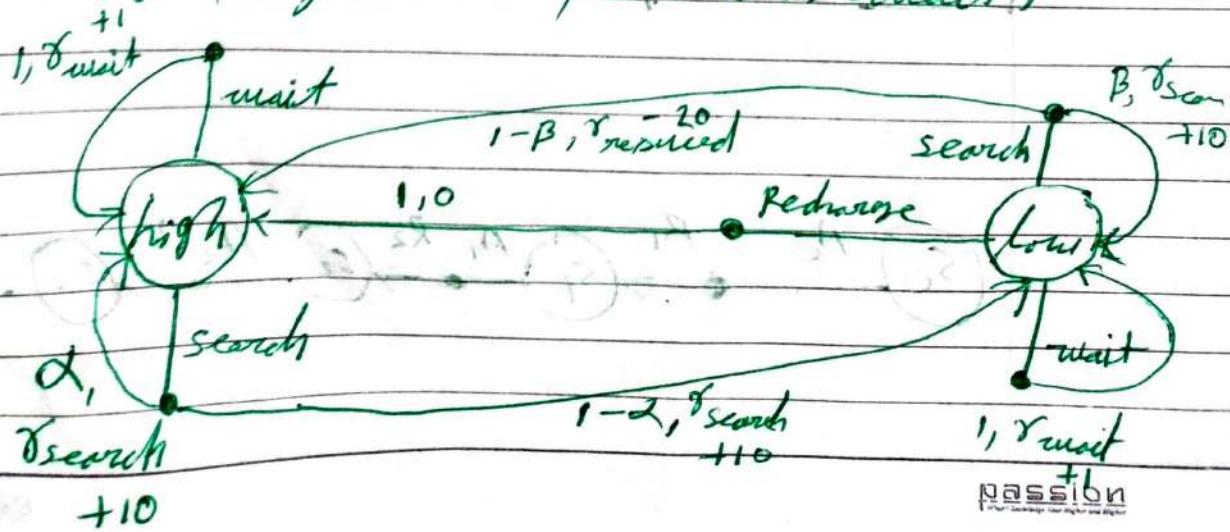
- \* The present state contains all the info necessary to predict the future

### Recycling Robot

$S_t = \{ \text{low}, \text{high} \}$

$A_t(\text{low}) = \{ \text{Search}, \text{Wait}, \text{Recharge} \}$

$A_t(\text{high}) = \{ \text{Search}, \text{Wait} \}$



Goal of an Agent:

maximize the expected return

$$E[G_t] = E[R_{t+1} + R_{t+2} + \dots + R_T]$$

↑  
final  
time  
step

- \* In episodic tasks, the agent-environment interaction breaks up into episodes
- \* Each episodes have terminal state

Reward hypothesis:

- \* Identify where reward signals come from
- \* develop algo that search the space of behaviors to maximize reward signals.

Episodic task  
interaction breaks naturally into episodes

\* each episode ends in a terminal state

\* episodes are independent

continuing task  
Interaction goes on continually

No terminal state

# MDP: sequential decision making

Date - 1 - 1 -

Continuing Tasks

$$G_t = R_{t+1} + R_{t+2} + \dots$$

Discounting

$$\begin{aligned} G_t &= R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{k-1} R_{t+k} + \dots \\ &= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \end{aligned}$$

How to make sure  $G_t$  is finite?

Discount the reward in the future by  $\gamma$  where

$$0 < \gamma < 1$$

Assume  $R_{\max}$  is the maximum reward the agent can receive

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \leq \sum_{k=0}^{\infty} \gamma^k R_{\max} =$$

$$R_{\max} \sum_{k=0}^{\infty} \gamma^k$$

Geometric series

$$= R_{\max} \times \frac{1}{1-\gamma} \rightarrow \text{finite}$$

Date - / - /

Effect of  $\gamma$  on agent behavior

$$\gamma = 0$$

$$G_t \doteq R_{t+1}$$

Agent only cares about the immediate reward  $\rightarrow$  short-sighted agent

$$\gamma \rightarrow 1$$

Agent take future rewards into account more strongly.

\* Rewardive nature of rewards:

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

$$= R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots)$$

\*  $G_t = R_{t+1} + \gamma G_{t+1}$

+ A policy maps each state to a single action. This kind of policy is called the deterministic policy.

$$\pi(s) = a$$

states      action

$s_0 \xrightarrow{\pi(a_0)} a_0$

$s_1 \xrightarrow{\pi(a_1)} a_1$

$s_2 \xrightarrow{\pi(a_2)} a_2$

state      Action

$s_0 \quad a_1$

$s_1 \quad a_0$

$s_2 \quad a_0$

Date 20/04/2021

- \* A stochastic policy is one where, multiple actions may be selected with non-zero probability.

$$\pi(a|s)$$

$$\sum_{a \in \mathcal{A}(s)} \pi(a|s) = 1$$

$$\pi(a|s) \geq 0$$

- \* policy can only depend on the current state  
(valid & invalid policies)

- - - - state value function - - - -

a state value function is the future reward an agent can expect to receive starting from a particular state, more precisely the state value function is the expected return from a given state.

$$V_\pi(s) \doteq \mathbb{E}_\pi [G_t | s_t = s]$$

Pi indicate the value function is contingent on the agent

Selecting actions according to  $\pi$ .

--- Action-value function ---

$$Q_{\pi}(s, a) \doteq E_{\pi}[G_t | S_t = s, A_t = a]$$

An action value describes what happens when the agent first selects a particular action.

$$\gamma = 0.9$$

	A	B	-1	3 3	8.8	5.3
-1	+10	+5	25%	1.3	2.0	2.1
		B	25% ↘ 25% ↗	0.1	0.7	-0.4
	A'		policy	-1.9	-1.3	-1.4

- \* State-value functions represent the expected return from a given state under a specific policy.
- \* Action value function represent the expected return from a given state after taking a specific action later following a specific policy.

--- Bellman Equation ---

The Bellman equation for the state value function defines a relationship between the value of a state and the value of his possible successor state.

- - - state-value Bellman equation - - -

$$V_n(s) \doteq E_n[(r_{t+1} | S_{t+1} = s)]$$

$$= E_n[R_{t+1} + \gamma V_n(S_{t+1}) | S_t = s]$$

- - -> first we expand the expected return as a sum over possible action choices made by the agent

- - -> second we expand over possible rewards and next states conditioned on state  $s$  and action  $a$ . Little  $\delta$  represents each possible reward outcome.

$$= \sum_a \pi(a|s) \sum_{s'} \sum_{\delta} P(s', r | s, a) [\delta + \gamma E_n[R_{t+1} | S_{t+1} = s']]$$

$$= \sum_a \pi(a|s) \sum_{s'} \sum_{\delta} P(s', r | s, a) [\delta + \gamma V_n(s')]$$

- - - Action-value Bellman Equation - - -

Value of a state-action pair in terms of its possible successors state-action pair.

In this case, the equation does not begin with the policy selecting an action because action

is already fixed as part of the state action pair. Instead we skip directly to the dynamic function  $p$  to select the immediate reward and next state  $s'$ .

$$q_{\pi}(s, a) \doteq E_{\pi}[G_t | S_t = s, A_t = a]$$

$$= \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma E_{\pi}[G_{t+1} | S_{t+1} = s']]$$

we can't stop here, we want to recursive equation for the values of one state action pair in terms of next state action pair.

At the moment, we have the expected return given only the next state. To change this, we can express the ~~expected~~ expected return from the next state as a sum of the agent's possible actions ~~choices~~ choices.

In particular, we can change the expectation to be conditioned on both the next state & the next action and then sum over all possible actions.

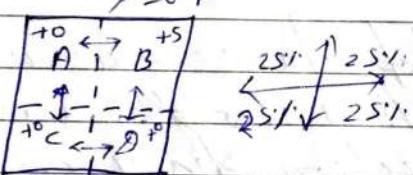
Bellman equation does not work with larger

Date - 1/1/178 P

$$\sum_{s'} \sum_{\pi} P(s', \pi | s, a) [\pi + \gamma \sum_{a'} \pi(a' | s') E_n(h_{n+1} | s_{n+1})]$$

$$= \sum_{s'} \sum_{\pi} P(s', \pi | s, a) [\pi + \gamma \sum_{a'} \pi(a' | s') h_n(s', a)]$$

Example : Gridworld



$$V_n(A) \doteq E_n(G_{n+1} | S_{n+1} = A)$$

$$V_n(B) \doteq E_n(G_{n+1} | S_{n+1} = B)$$

$$V_n(A) = \sum_a \pi(a | A) (\pi + 0.7 V_n(s'))$$

(for each action there's only one possible associated next state and reward.  
That's the sum over  $s'$  and  $\pi$  reduces to a single value.

$$V_n(A) = \sum_a \pi(a | A) (\pi + 0.7 V_n(s'))$$

$$V_n(A) = \frac{1}{4} (5 + 0.7 V_n(B) + \frac{1}{4} 0.7 V_n(C) + \frac{1}{2} 0.7 V_n(D))$$

$$V_n(B) = \frac{1}{2} (5 + 0.7 V_n(B) + \frac{1}{4} 0.7 V_n(D) + \frac{1}{4} 0.7 V_n(S))$$

$$= s', A_{t+1} = a' ]$$

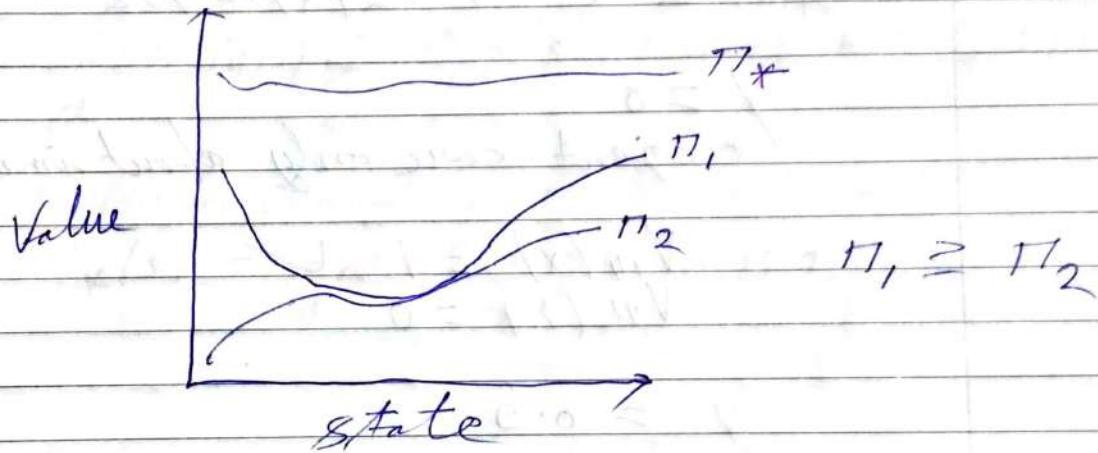
$$V_n(A) = 4.2$$

$$V_n(B) = 6.1$$

$$V_n(C) = 2.2$$

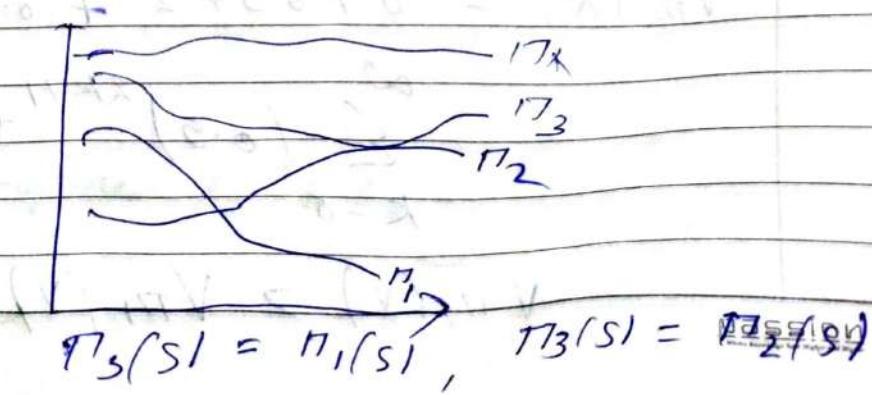
$$V_n(D) = 4.2$$

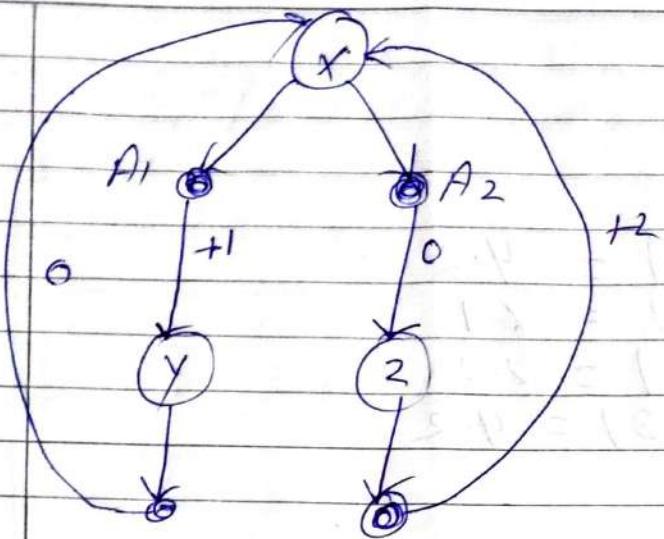
\* optimal policies



\* An optimal policy  $\pi_*$  is as good as or better than all other policies.

\* There is always an optimal policy





$$\pi_1(x) = A_1 \quad \pi_2(x) = A_2$$

$$\gamma = 0$$

\* agent care only about immediate reward

$$V_{\pi_1}(x) = 1 \quad \checkmark$$

$$V_{\pi_2}(x) = 0$$

$$\gamma = 0.9$$

$$V_{\pi_1}(x) = 1 + 0.9 * 0 + (0.9)^2 * 1 + \dots$$

$$\sum_{k=0}^{\infty} (0.9)^{2k} = \frac{1}{1 - 0.9^2} \approx 5.3$$

$$V_{\pi_2}(x) = 0 + 0.9 * 2 + (0.9)^2 * 0 + \dots$$

$$\sum_{k=0}^{\infty} (0.9)^{2k+1} * 2 \approx 9.5$$

$$V_{\pi_2}(x) > V_{\pi_1}(x)$$

\* An optimal policy is the policy with the highest possible value function in all states.

$$\forall \pi^* \quad V_{\pi^*}(s) = E_{\pi^*} [G_t | S_t = s] = \max_a V_{\pi^*}(s)$$

$$\forall \pi^* \quad Q_{\pi^*}(s, a) = \max_{\pi} Q_{\pi}(s, a) \text{ for all } s \in S \text{ and } a \in A$$

$$V^*(s) = \sum_a \pi^*(a|s) \sum_{s'} \sum_{\pi} p(s', r | s, a) [r + \gamma V^*(s')]$$

Bellman Optimality equation for  $V^*$  &  $Q^*$

$$V^*(s) = \max_a \sum_{s'} \sum_{\pi} p(s', r | s, a) [r + \gamma V^*(s')]$$

$$Q^*(s, a) = \sum_{s'} \sum_{\pi} p(s', r | s, a) [r + \gamma \max_a Q^*(s', a)]$$

optimal policy

$$\pi^*(s) = \arg \max_a \sum_{s'} \sum_{\pi} p(s', r | s, a) [r + \gamma V^*(s')]$$

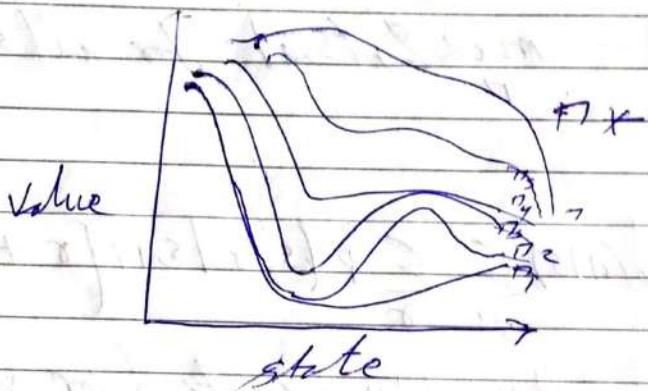
$$\pi^*(s) = \arg \max_a q^*(s, a)$$

week - 04

0

Date \_\_\_\_\_  
Date - 1 - 1 -

- \* policy evaluation is the task of determining state-value function  $V_{\pi}$ , for a particular policy  $\pi$
- \* Control is the task of improving an existing policy.



$\pi, P, r \rightarrow [DP] \rightarrow V_{\pi}$  policy evaluation

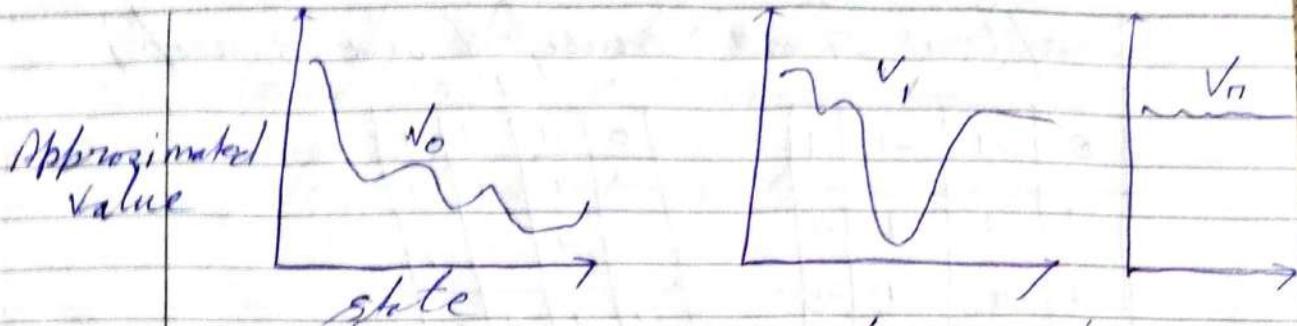
$P, r \rightarrow [DP] \rightarrow \pi^*$  control

--- Iterative policy evaluation ---

use bellman equation as update rule

$$V_{k+1}(s) \leftarrow \sum_a \pi(a|s) \sum_{s'} P(s'|r|s,a) [r + \gamma V_k(s')]$$

(sequence of better estimation for value function)

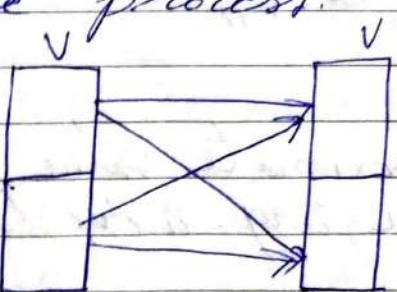


for any  $\lambda_0$

$$\lim_{k \rightarrow \infty} v_k = v_n$$

We store two arrays, each has one entry for every state. One array which we label  $V$  stores the current approximation value function. Another array  $V'$  stores the updated value.

By using 2 arrays, we can compute new values from the old one state at a time without the old values being changed in the process.



0	0	0	0
0	0	0	0
0	0	0	0
0	0	0	0

```

graph LR
    start(( )) --> q1
    q1 -- "0" --> q1
    q1 -- "1" --> q2
    q2 -- "1" --> q1
    q2 -- "1" --> q2
    style start fill:none,stroke:none
    style q1 fill:none,stroke:none
    style q2 fill:none,stroke:none
    
```

$$0.25 \times (-1+0) + 0.25 \times (-1+0) + 0.25 \times (-1+0) + 0.25 \times (-1+0) = -1$$

Date - / - / -

after first sweep (one sweep)

0	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-10	0

0	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	-1
-1	-1	-1	0

\* for multiple sweeps algorithm

[Input  $\pi$ , the policy to be evaluated]

$$V \leftarrow \vec{0}, V' \leftarrow \vec{0}$$

Loop:

$\Delta \leftarrow 0$   
Loop for each  $s \in S$ :

$$V'(s) \leftarrow \text{(formula written on previous page)}$$
$$\Delta \leftarrow \max(\Delta, |V'(s) - V(s)|)$$

$$V \leftarrow V'$$

until  $\Delta < \phi$  (a small +ve number)  
Output  $V \approx V_\pi$

once the approximate value function stops changing, we've converge to  $V_\pi$

\* policy improvement

$$\pi(s) = \arg\max_a \sum_{s' r} p(s'|r|s,a)[r + \gamma V_\pi(s')] \text{ for all } s \in S$$

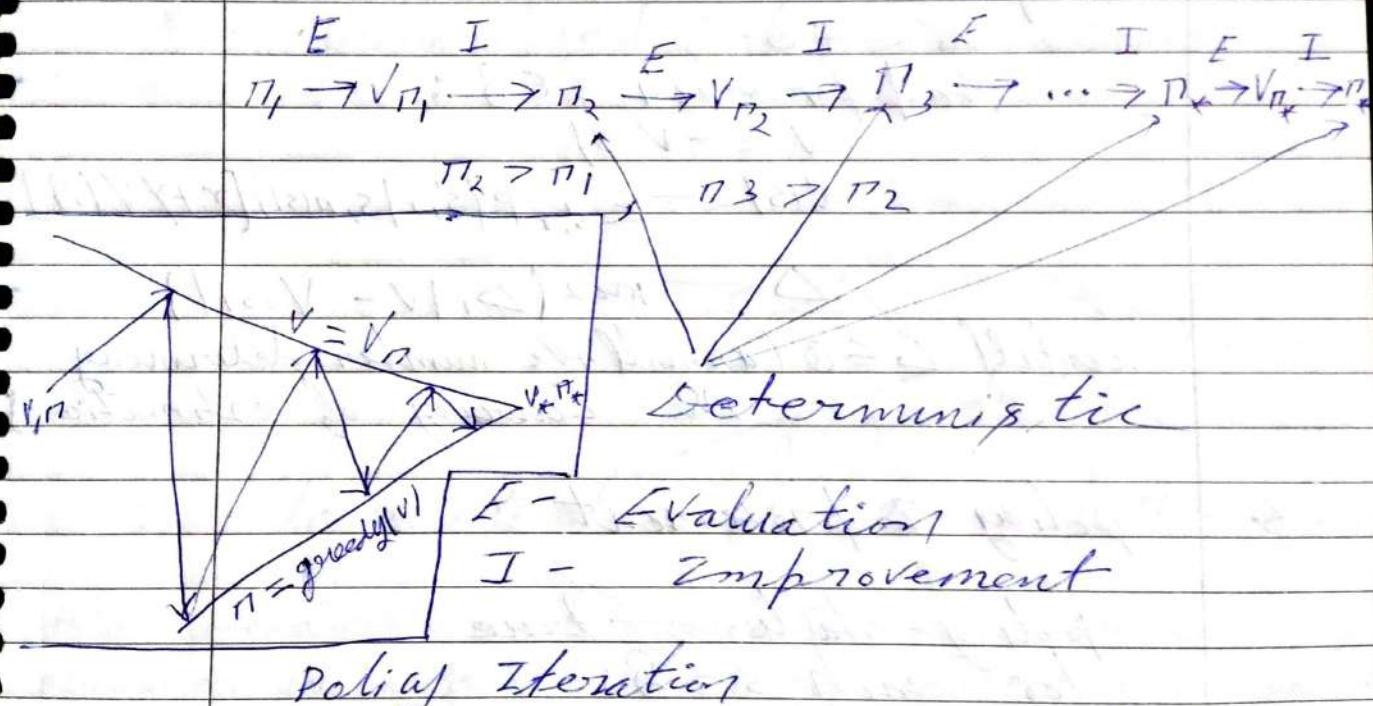
Date - 1 - 1 -

$\rightarrow V_{\pi}$  obeys the bellman optimality equation  $\rightarrow \pi$  is optimal

policy Improvement Theorem

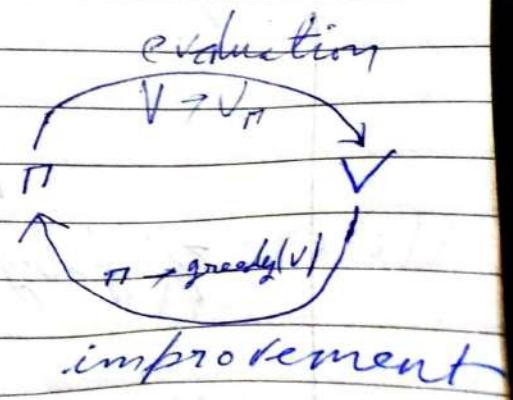
$$q_{\pi}(s, \pi'(s)) \geq q_{\pi}(s, \pi(s)) \text{ for all } s \in S$$

\* The new policy is a strict improvement over  $\pi$  unless  $\pi$  is already optimal.



Policy Iteration

$$\begin{aligned} \pi_2 &\leftarrow \cancel{V_{\pi_1}} \\ \pi_2 &\leftarrow V_{\pi_1} \\ \pi_2 &\rightarrow V_{\pi_2} \\ \pi_k &\leftrightarrow V_k \end{aligned}$$



Date - 1 - 1 -

Policy Iteration (using iterative policy evaluation) for estimating  $\pi \approx \pi^*$

### 1. Initialization

$V(s) \in R$  and  $\pi(s) \in A(s)$  arbitrarily for all  $s \in S$

### 2. policy Evaluation

Loop:

$$\Delta \leftarrow 0$$

loop for each  $s \in S$ :

$$V' \leftarrow V(s)$$

$$V(s) \leftarrow \sum_{s', r} P(s', r | s, \pi(s)) [r + \gamma V'(s')]$$

$$\Delta \leftarrow \max(\Delta, |V - V'|)$$

until  $\Delta < 0$  (a small +ve number determining the accuracy of estimation)

### 3. policy Improvement

policy-stable  $\leftarrow$  true  
for each  $s \in S$ :

$$\text{old-action} \leftarrow \pi(s)$$

$$\pi(s) \leftarrow \arg\max_a \sum_{s', r} P(s', r | s, a) [r + \gamma V(s')]$$

if  $\text{old-action} \neq \pi(s)$ , then policy-stable  $\leftarrow$  false

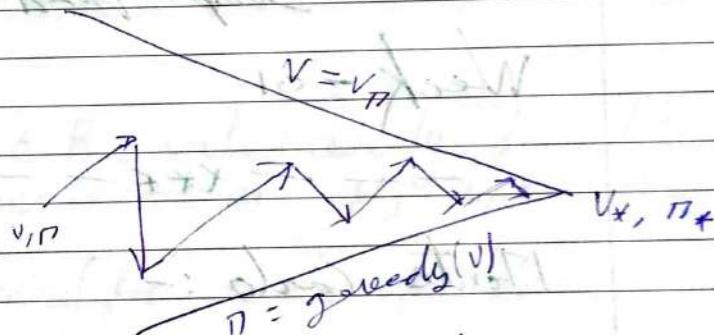
if policy-stable, then stop and return  $V \approx V^*$   
and  $\pi \approx \pi^*$ ;  
else go to 2

Date - 1 - 1 -

\* Policy iteration works by alternating policy evaluation and policy improvement

\* Policy iteration follows a sequence of better and better policies and value functions until it reaches the optimal policy and associated optimal value functions.

\* Generalized policy iteration



\* Value iteration, for estimating  $\pi \approx \pi^*$

Alg. parameter: a small threshold  $\delta > 0$  determining accuracy of estimation  
initialise  $V(s)$ , for all  $s \in S^*$ , except that  $V(\text{terminal}) = 0$   
loop:

$$\Delta \leftarrow 0$$

loop for each  $s \in S^*$ :

$$V \leftarrow V(s)$$

$$V(s) \leftarrow \max_a \sum s' r P(s, a)[\gamma + V(s')]$$

$$\Delta \leftarrow \max(\Delta, |V - V(s)|)$$

until  $\Delta < \delta$

output a deterministic policy  $\pi \approx \pi^*$ , such

that

$$\pi(s) = \arg \max_a \sum s' r P(s, a)[\gamma + V(s')]$$

Noah

Date - 1 - 1 -

## Summary

- \* Value Iteration allows us to combine policy evaluation and improvement into a single update.
- \* Asynchronous dynamic programming method give us the freedom to update states in any order

## Course - 2

### Sample based learning methods

#### Week - 01

— — — X \*\* — — —

#### Monte Carlo :-

Monte Carlo methods estimate value by averaging over a large number of random samples

So, a monte carlo method for learning a value function would first observe multiple returns from the same state.

Then it averages those observed returns to estimate the expected return from that state.

As the number of samples

increases, the average tends to get closer and closer to the expected return. The more returns the agent observe from a state, the more likely it is that the sample average is close to the state value.

MC prediction, for estimating  $V = V_{\pi}$

Input: a policy  $\pi$  to be evaluated

Initialize:

$V(s) \in \mathbb{R}$ , arbitrarily, for all  $s \in S$   
 $\text{Returns}(s) \leftarrow$  an empty list, "

Loop forever (for each episode):

Generate an episode following  $\pi$ :  $S_0, A_0, R_1, \dots$

$G_t \leftarrow 0$

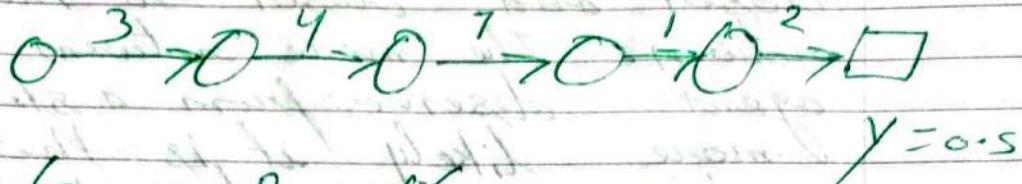
Loop for each step of episode,  $t = T-1, T-2, \dots$

(a)  $G_t \leftarrow G_t + R_{t+1}$

Append  $G_t$  to Returns( $S_t$ )

$V(S_t) \leftarrow \text{average}(\text{Returns}(S_t))$

## Computing returns efficiently



$$G_0 = R_1 + y G_1$$

$$G_1 = R_2 + y G_2$$

$$G_2 = R_3 + y G_3$$

$$G_3 = R_4 + y G_4$$

$$G_4 = R_5 + y G_5$$

$$G_5 = 0$$

\* Compute backward to save computation

$$G_5 = 0$$

$$G_4 = R_5 + 0.5 \cdot 0 = 2$$

$$G_3 = R_4 + y R_5 = 2$$

$$G_2 = R_3 + y R_4 + y^2 R_5 = 8$$

$$G_1 = R_2 + y R_3 + y^2 R_4 + y^3 R_5 = 8$$

$$G_0 = R_1 + y R_2 + y^2 R_3 + y^3 R_4 + y^4 R_5 = 7$$

(Table Acc, sum, order)

Returns (s)

V(s)

A

Returns (A) = [1]

1

B

Returns (B) = [1]

1

CP

$R = -1, 0, +1$

Date: 1/1/1

Dealer |

10

7  6

player |  13  $G_{10} = 1$

Dealer |

10

1  6  7

player |  20  $G_1 = 1$

Dealer |  14

10  4

7  10  7

player |  20

Date - / - / -

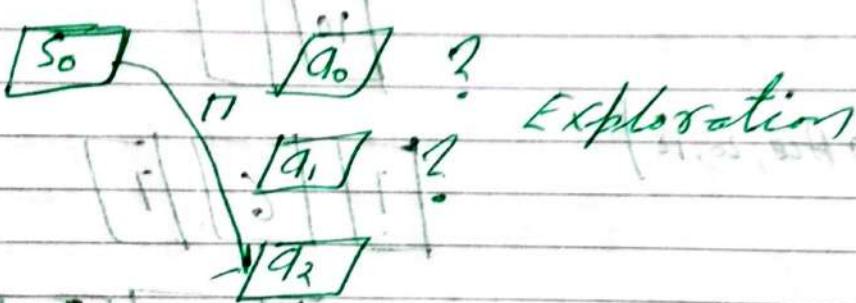
[Dealer] [23]

[10] [4] [9]

[7] [6] [1]

[Player] [20]  $R = 1$

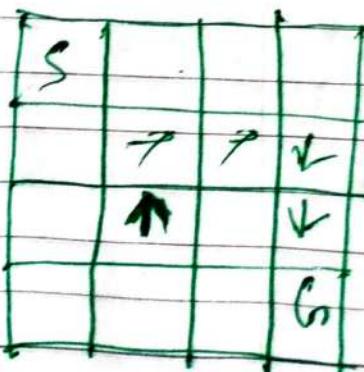
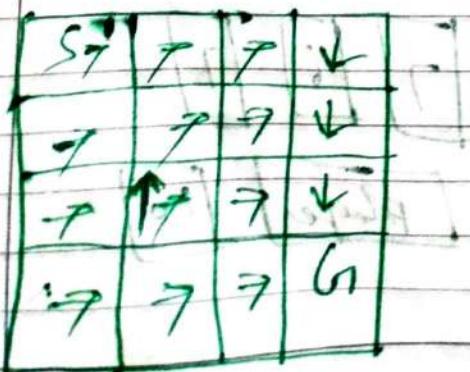
## \* Monte-Carlo Action-Values



## \* Exploring Starts

$S_0, A_0, S_1, S_2 \dots$

Random Explore and P



Date - / - / -

## Monte Carlo Generalized Policy Iteration

$$\pi_0 \rightarrow \pi_1 \rightarrow \pi_2 \rightarrow \dots$$

Improvement:

$$\pi_{k+1}(s) = \underset{a}{\operatorname{argmax}} Q_k(s, a)$$

Evaluation:

Monte Carlo Prediction

Monte Carlo ES (Exploring Starts), for  $\pi \leftarrow \pi^*$

Initialize:

$$\begin{aligned} \pi(s) &\in A(s), \text{ for all } s \in S \\ Q(s, a) &\in R, \text{ for all } s \in S, a \in A(s) \\ \text{Returns}(s, a) &\leftarrow \text{empty list}, \quad \Pi \end{aligned}$$

Loop forever (for each episode)

choose  $s \in S, a \in A(s)$ , randomly such that  
all pairs have probability  $> 0$

Generate an episode from  $s_0, a_0$ , following  $\pi(s_0, a_0)$

$$G_t \leftarrow 0$$

Loop for each step of episode,  $t = T-1, T-2, \dots$

$$G_t \leftarrow \gamma G_t + R_{t+1}$$

Append  $G_t$  to Returns( $s_t, a_t$ )

$$Q(s_t, a_t) \leftarrow \underset{\pi}{\operatorname{average}}(\text{Returns}(s_t, a_t))$$

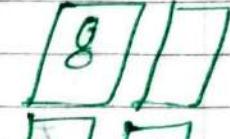
$$\pi(s_t) \leftarrow \underset{a}{\operatorname{argmax}} Q(s_t, a)$$

Date - / - / -

$s, a$	>Returns( $s, a$ )	$Q(s, a)$	$r(s)$	Dealer	<input type="checkbox"/>
		Hit / Stick			

$x, \text{stick}$  Ret( $x, \text{stick}$ ) = 0 1 Stick  
 $y, \text{hit}$  Ret( $y, \text{hit}$ ) = 1 0 Hit

$$y = (\text{None}, 13, 8)$$



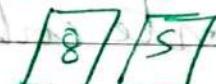
Player [13]  $R_{10} =$

Dealer



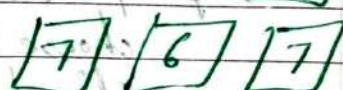
Player [20]  $R_{11} =$

Dealer [13]



Player [20]  $\circlearrowleft$

Dealer [22]



Player [20]  $R = 1$

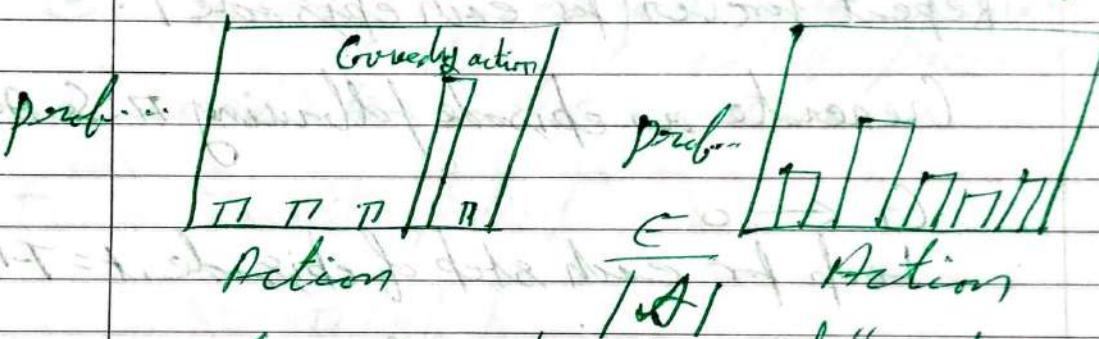
Date - / - / -

## \* Epsilon - Soft policies

$\epsilon$  greedy policies are a subset of a larger class of policies called  $\epsilon$  soft policies.  $\epsilon$  soft policies take each action with prob at least  $\epsilon$  over the number of actions.

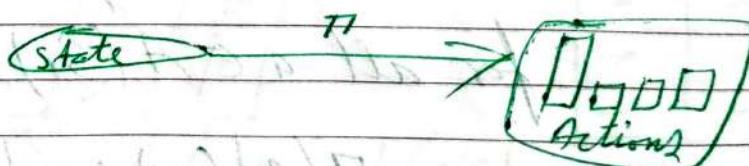
### $\epsilon$ -Greedy Exploration

$\epsilon$ -Greedy policies  $\subset$   $\epsilon$ -soft policies



Thus we don't need "Exploring starts"

$\epsilon$ -soft policies are always stochastic



$\epsilon$ -greedy policies and deterministic policies  
deterministic

Sr	$\rightarrow$	$\uparrow$	$\downarrow$
$\downarrow$	$\leftarrow$	$\leftarrow$	$\leftarrow$
C1			

Sr	$\rightarrow$	$\leftarrow$	$\uparrow$	$\downarrow$
$\leftarrow$	$\leftrightarrow$	$\leftrightarrow$	$\leftrightarrow$	$\leftrightarrow$
C1				

$\epsilon$ -greedy

PASSION

Date - 1-1-

\*  $\epsilon$ -soft policies may not be optimal

UK control (for  $\epsilon$ -soft policies),  $\pi = \pi_\epsilon$

algo parameters: small  $\epsilon > 0$   
initialize

$\pi \leftarrow \epsilon$ -soft policy  
 $Q(s,a) \in \mathbb{R}$ , for all  $s \in S, a \in A(s)$   
Returns( $s, a$ )  $\leftarrow$  empty list

Repeat forever (for each episode):

Generate an episode following  $\pi$ :  $S_0, A_0, \dots$

$G_t \leftarrow 0$

Loop for each step of episode,  $t = T-1, T-2, \dots$

$G_t \leftarrow \gamma G_t + R_{t+1}$

Append  $G_t$  to Returns( $S_t, A_t$ )

$Q(S_t, A_t) \leftarrow \text{average}(\text{Returns}(S_t, A_t))$

$A^* \leftarrow \underset{a}{\operatorname{argmax}} Q(S_t, a)$

for all  $a \in A(S_t)$ :

$\pi(a|S_t) \leftarrow (1-\epsilon) + \epsilon / |A(S_t)|$   
 $\epsilon / |A(S_t)|$

$\pi(a = A^*)$

$\pi(a \neq A^*)$

## $\epsilon$ -soft policy

- Values based on suboptimal policy
- Actions " " " "

## \* On-Policy & off-policy

- On-Policy : improve and evaluate the policy being used to select actions.
- Off-policy : improve and evaluate a different policy from the one used to select actions.

## \* Target policy

$\pi_{\text{tgt}}$

- Learn values for this policy
- for example the optimal " "

↑	↑	↑	•
↑	↑	↑	↓
↑	↓	↓	↑

↑	↑	•
↑	↑	↑
↑	↑	↑
↑	↑	↑

## \* Behavior policy $\pi_b$

- Select actions from this policy
- Generally an exploratory policy (uniform random)

\* One key rule of off policy learning is that the behavior policy must cover the target policy.

$$\pi(a|s) > 0 \text{ where } b(a|s) > 0$$

9	7	0
↑		↑
↓	→	↑



Consider this state with the behavior policy always goes up but the target policy goes to right agent can not learn the correct action value for that state b/c never observed samples of what would happen if it goes right

~~on off~~

\* on policy

		0
↑	↑	↑↑
↑	↑	↑↑

$$\pi(a|s) = b(a|s)$$

- \* Importance sampling uses samples from one probability distribution to estimate the expectation of a different distribution.
- \* Derivation of importance sampling

Sample: n numb  
Estimates:  $E_p[X]$

$$E_p[X] = \underbrace{\sum_{n \in X} n p(n)}_{\text{Sample}}$$

$$= \sum_{n \in X} n p(n) \frac{b(n)}{b(n)}$$

$$= \underbrace{\sum_{n \in X} n p(n) b(n)}_{\text{Importance Sampling}}$$

$p(n) = \frac{n}{b(n)}$  importance sampling ratio

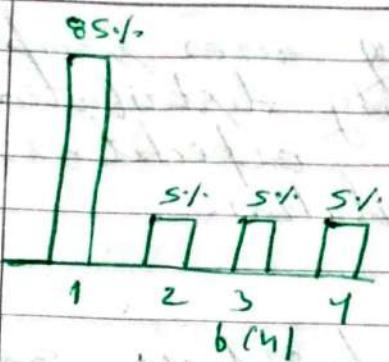
$$= E_b[X P(X)]$$

$$E[X] \approx \frac{1}{n} \sum_{i=1}^n n_i$$

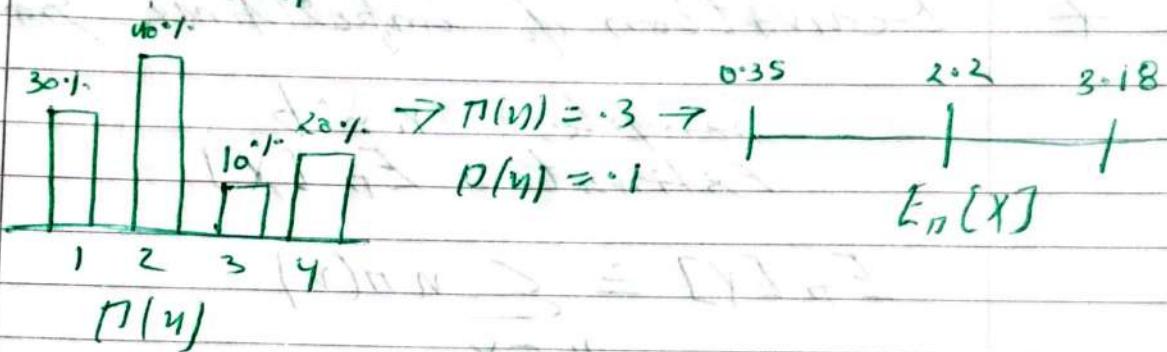
$$\approx \frac{1}{n} \sum_{i=1}^n n_i p(n_i)$$

n, numb

Date - / - / -



$$\begin{aligned} n &= 1 \rightarrow x = [1] \rightarrow x = [1, 3] \rightarrow \\ b(\eta) &\approx 85 \\ n &= 3 \\ b(\eta) &\approx .05 \end{aligned}$$



$$\begin{aligned} \pi(\eta) &= .3 \rightarrow \\ p(\eta) &= .1 \\ 0.35 & \quad 2.2 \quad 3.18 \\ E_n(x) & \end{aligned}$$

$$\frac{1}{n} \sum_{i=1}^n \pi(\eta_i) \rightarrow 1x \frac{.3}{.85} = 0.35$$

$$\frac{(1x \frac{.3}{.85}) + (3x \frac{.1}{.05})}{2} = 3.18$$

$$\frac{(1x \frac{.3}{.85}) + (3x \frac{.1}{.05}) + (1x \frac{.3}{.85})}{3} = 2.24$$

$E(X) \approx 1.2$

\* off-policy Monte Carlo

$$P = \frac{P(\text{trajectory under } \pi)}{P(\text{trajectory under } \delta)}$$

$$V_{\pi}(s) = E_b [PG_t | S_t = s]$$

off-policy Trajectories

$$\rightarrow P(A_t, S_{t+1}, A_{t+1}, S_{t+2}, \dots, S_T | S_t, A_{1:T}) =$$

$$\rightarrow b(A_t | S_t) P(S_{t+1} | S_t, A_t)$$

$$b(A_{t+1} | S_{t+1}) P(S_{t+2} | S_{t+1}, A_{t+1}) \dots$$

$$\cancel{P(S_T | S_{T-1}, A_{T-1})}$$



T-1

$$\prod_{k=t}^{T-1} b(A_k | S_k) P(S_{k+1} | S_k, A_k)$$

K=t

P(trajectory number b)

$$\rho_{t:T-1} \doteq \frac{P(\text{traj under } \pi)}{P(\text{traj under } b)}$$

$$\doteq \prod_{k=t}^{T-1} \frac{P(A_k | S_k) P(S_{k+1} | S_k, A_k)}{b(A_k | S_k) P(S_{k+1} | S_k, A_k)}$$

$$P_{t:T-1} \doteq \prod_{k=t}^{T-1} \frac{\pi(A_k | S_k)}{b(A_k | S_k)}$$

$$\mathbb{E}_\theta [P_{t:T-1}(G_t) | S_t = s] = V_\pi(s)$$

off-policy MC prediction, for  $V \approx V_\pi$

Input: a policy  $\pi$  to be evaluated

Initialize:

$V(s) \in \mathbb{R}$ , arbitrarily, for all  $s \in S$

Returns<sub>t</sub>(s)  $\leftarrow$  an empty list,  $\{\}$

Loop forever (for each episode):

(generate an episode following b(S<sub>0</sub>, A<sub>0</sub>, R<sub>1</sub>, S<sub>1</sub>, ...))

( $t < 0$ )  $W \leftarrow 1$

(1) Loop for each step of episode,  $t=T-1, \dots, 0$

$$G_t \leftarrow \gamma W G_t + R_{t+1}$$

Append  $G_t$  to Returns<sub>t</sub>(S<sub>t</sub>)

$$V(S_t) \leftarrow \text{average}(\text{Returns}_t(S_t))$$

$$W_t \leftarrow W \frac{\pi(A_t | S_t)}{b(A_t | S_t)}$$

Computing  $P_{t:T-1}$  incrementally

$$P_{t:T-1} = P_t P_{t+1} P_{t+2} \dots P_{T-2} P_{T-1}$$

$$W_1 \leftarrow P_{t-1}$$

$$W_2 \leftarrow P_{t-1} P_{t-2}$$

$$W_3 \leftarrow P_{t-1} P_{t-2} P_{t-3}$$

$$W_{t+1} \leftarrow W_t P_t$$

\* Temporal Difference  $\rightarrow$

$$V(S_t) \leftarrow V(S_t) + \alpha [G_t - V(S_{t+1})]$$

TD can update its value estimates on each step of the episode. It does not have to wait for the episode to complete. It just has to remember the previous state.

$$V_\pi(s) \doteq E_\pi[G_t | S_t = s]$$

$$= E_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s]$$

$$V_\pi(s) = R_{t+1} + \gamma V_\pi(S_{t+1})$$

$$V(S_t) \leftarrow V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

Tabular TD(0) for estimating  $V_n$

Input: the policy  $\pi$  to be evaluated

algo params: step size  $\alpha \in (0, 1]$

Initialize  $V(s)$ , for all  $s \in S^+$ , arbitrarily except that  $V(\text{terminal}) = 0$

Loop for each episode:

Initialize  $s$

Loop for each step of episode

$A \leftarrow$  action given by  $\pi$  for  $s$

Take action  $A$ , observe  $R, s'$

$$V(s) \leftarrow V(s) + \alpha[R + \gamma V(s') - V(s)]$$

$s \leftarrow s'$

until  $s$  is terminal

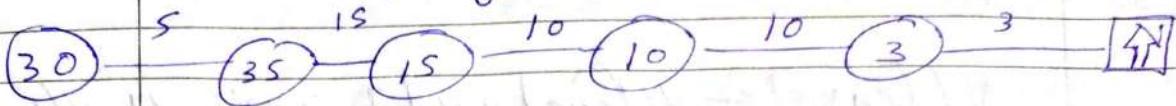
\* Temporal Difference learning - a way

to incrementally estimate the return through bootstrapping.

\* TD-error:  $\delta_t = R_{t+1} + \gamma V(s_{t+1}) - V(s_t)$

Date - 1 - 1 -

leave exit exit secondary  
highway road enter home arrive  
home street home



$$\alpha = 1$$

$$\gamma = 1$$

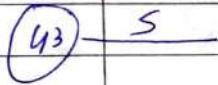
$$G_0 = 5 + 15 + 10 + 10 + 3 = 43$$

$$V(\text{leave}) \leftarrow V(\text{leave}) + \alpha [G_0 - V(\text{leave})]$$

30

43

30



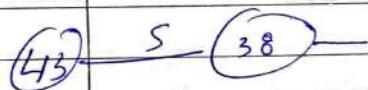
$$G_1 = 15 + 10 + 10 + 3 = 38$$

$$V(\text{exit}) \leftarrow V(\text{exit}) + \alpha [G_1 - V(\text{exit})]$$

35

38

35



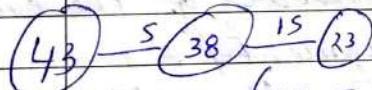
$$G_2 = 10 + 10 + 3 = 23$$

$$V(\text{exit highway}) \leftarrow V(\text{exit highway}) + \alpha [G_2 - V(\text{exit highway})]$$

15

23

15



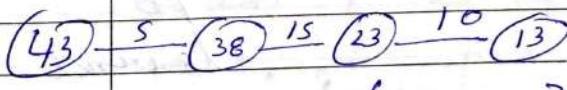
$$G_3 = 10 + 3 = 13$$

$$V(\text{secondary road}) \leftarrow V(\text{secondary road}) + \alpha [G_3 - V(\text{secondary road})]$$

10

13

10



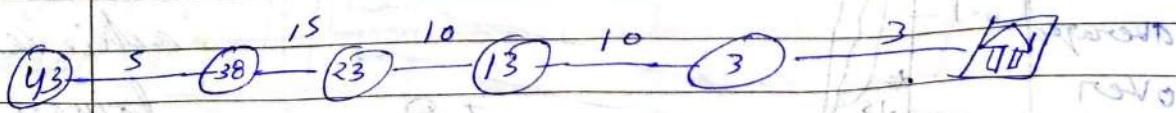
$$G_4 = 3$$

$$V(\text{enter home st}) \leftarrow V(\text{enter home st}) + \alpha [G_4 - V(\text{enter home st})]$$

3

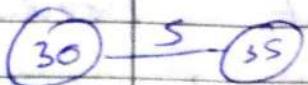
3

3

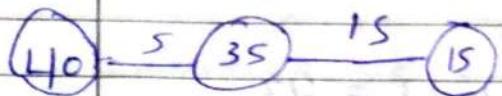


must wait until the end of the episode before learning can begin.

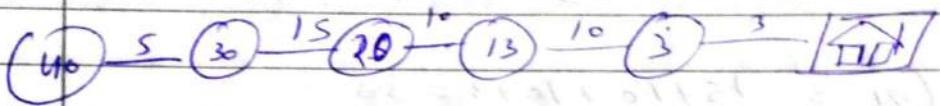
leave      exit



$$V(\text{leave}) \leftarrow V(\text{leave}) + \alpha [R_1 + \gamma V(\text{exit}) - V(\text{leave})]$$



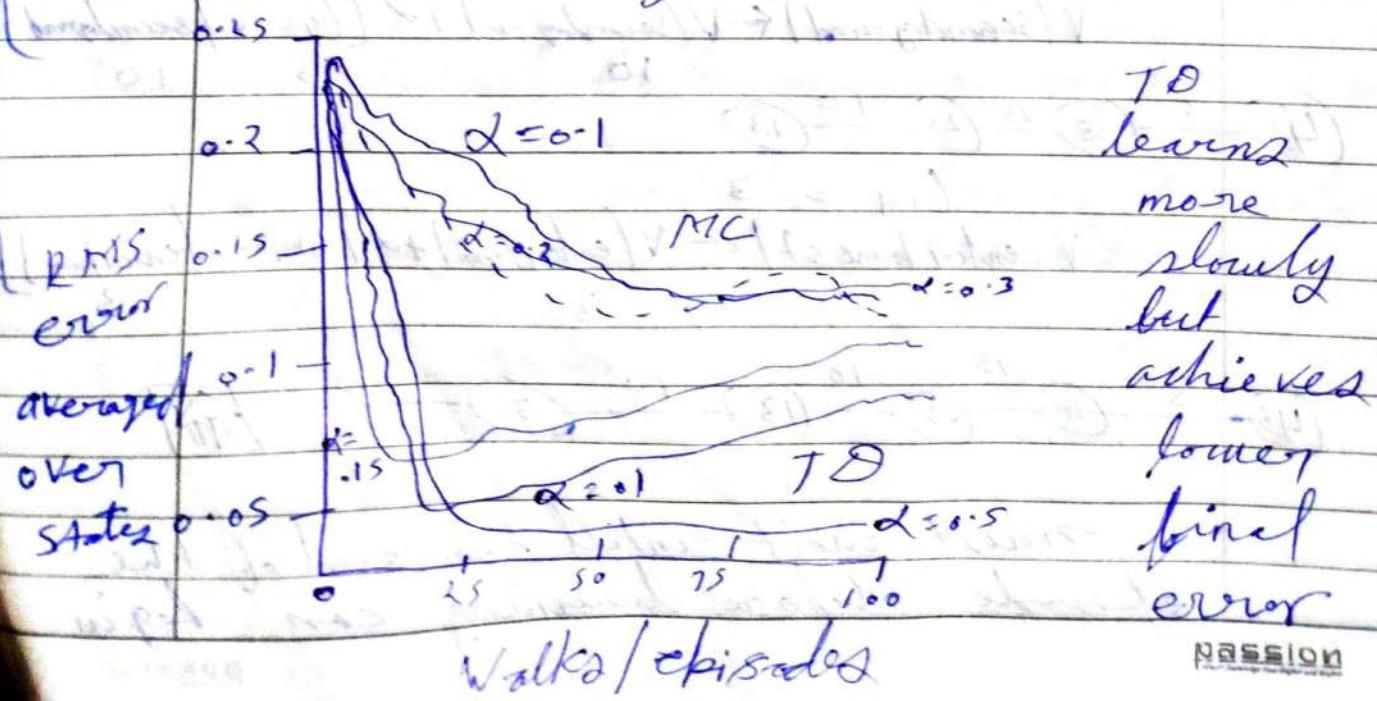
$$V(\text{exit}) \leftarrow V(\text{exit}) + \alpha [R_2 + \gamma V(\text{done}) - V(\text{exit})]$$



We can learn online without waiting for the episode to end.

## \* Advantages of TD methods

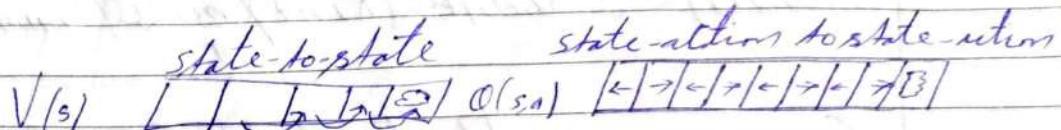
- (1) Don't require a model of the environment
- (2) Online and incremental.
- (3) Converge faster than MC Methods



Date - / - / -

SARSA :-

## From state-values to action-values



## The Sarrus algorithm

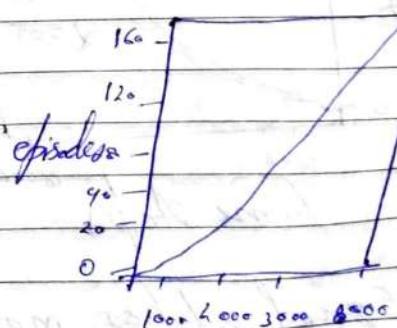
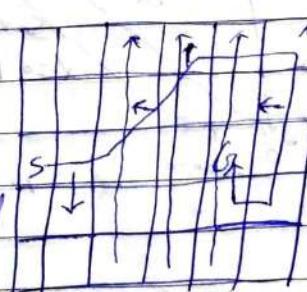
$S_t P_t R_{t+1} S_{t+1} A_{t+1}$

$$O(S_{t+1}, A_{t+1}) \leftarrow O(S_t, A_t) + \alpha \left( R_{t+1} + \gamma O(S_{t+1}, A_{t+1}) - O(S_t, A_t) \right)$$

In SARSA, the agent needs to know its next state action pair before updating its value estimates. That means it has to commit to its next action before the update.

An ~~of~~ epsilon of 0.1 means we take a random action 1 in every 10 steps.

# The windy windmill



Sarsa :

$$\underline{E \text{ o.1}}$$

205

## $\alpha$ -Learning

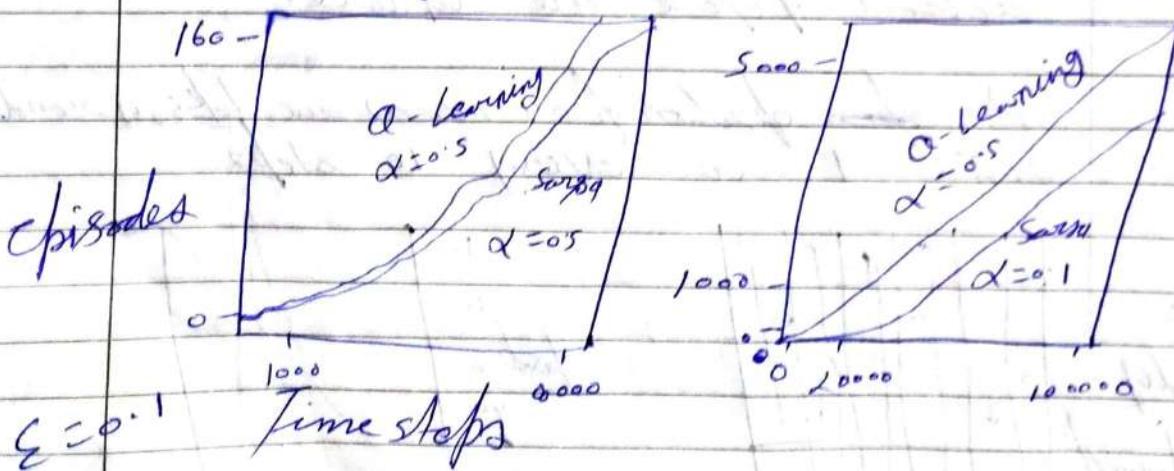
$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (R_{t+1} + \gamma \max_{A'} Q(S_{t+1}, A')) - Q(S_t, A_t)$$

Sarsa is sample-based version of policy iteration which uses Bellman equations for action values, that each depend on a fixed policy.

$\alpha$  Learning is a ~~sampled~~ sample-based version of value iteration which iteratively applies the Bellman optimality equation.

## Sarsa - Policy Iteration

## $\alpha$ Learning - Value Iteration

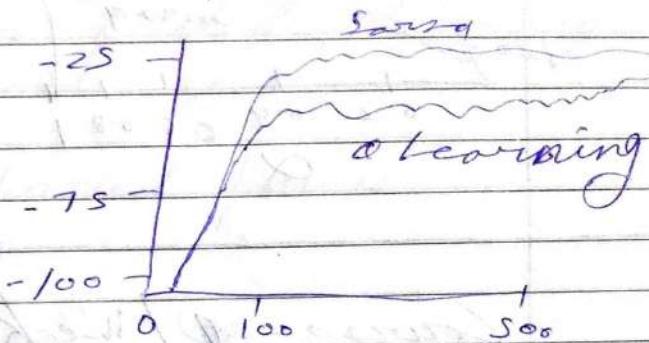


Equal slopes mean that both agents are completing episodes at the same time.

Sarsa learns same policy as  $\alpha$ -learning, but more slowly

Date - 1-1-

## The cliff walking environment



7) ~~α leaning~~ learning is off-policy without using importance sampling.

7) learning on-policy or off-policy may perform differently in control.

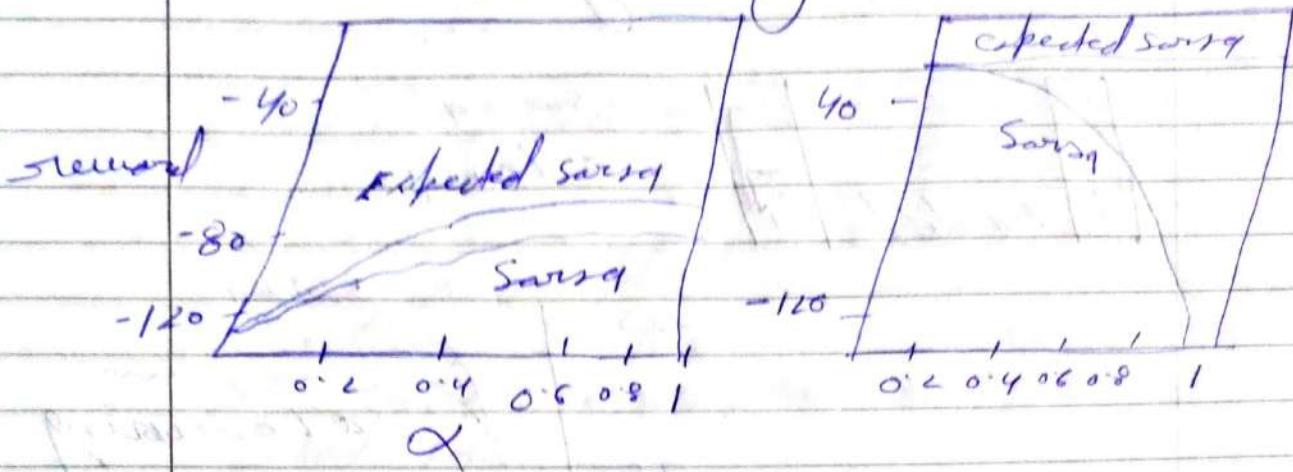
Expected Sarsa algo-

$$\hat{Q}(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (R_{t+1} + \gamma \sum_a \pi(a|S_{t+1}) \hat{Q}(S_{t+1}, a) - Q(S_t, A_t))$$

The algo is nearly identical to Sarsa, except the  $\hat{Q}$  error uses the expected estimate of the next action value instead of a sample of the next action value. That means that on every time step, the agent has to average the next state's action values according to how likely ~~they are~~ under the policy

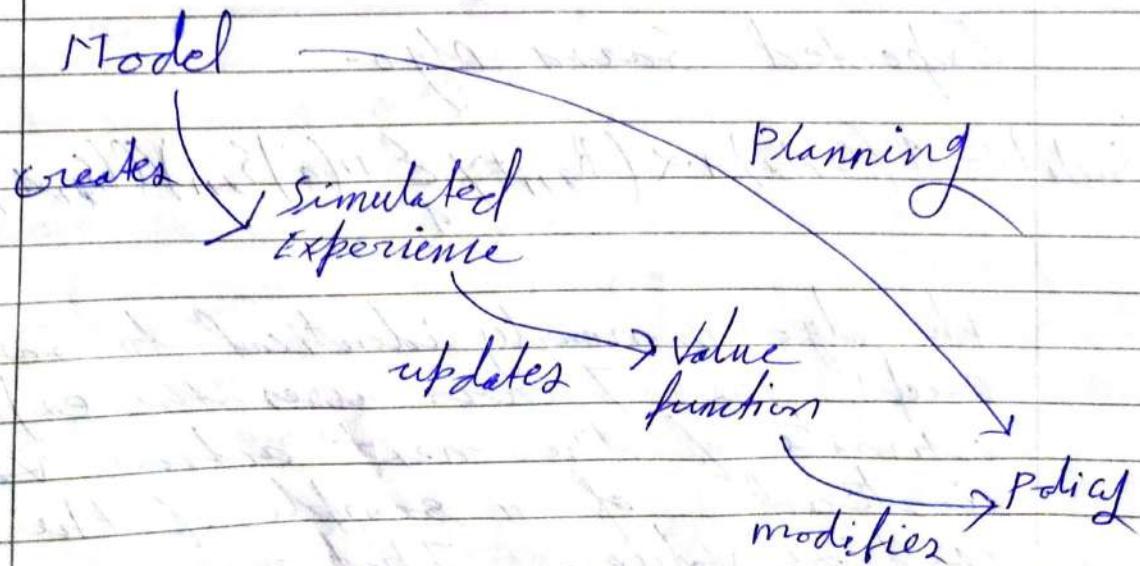
Date - / - / -

## The cliff-walking environment



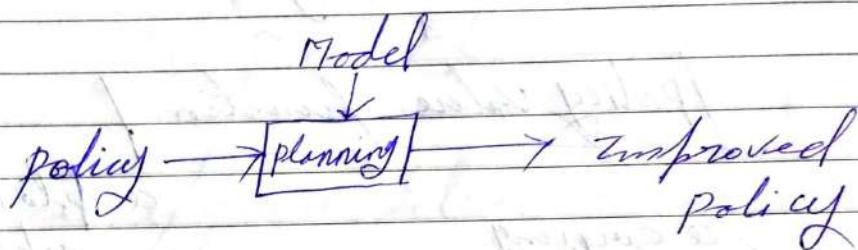
Course - 43 / Week - 4

Planning refers to the process of using a model to improve a policy.

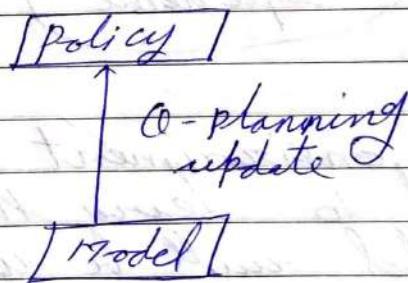


- \* Sample models procedurally generate samples, without explicitly storing the probability of each outcome
- \* Distribution models, contain a list of all outcomes and their probabilities thus require a lot of memory but gives good result.

planning improved policies



connection with Q-learning



Random-Sample one-step tabular Q-planning

1. Sampling  $S \rightarrow s \rightarrow$  [model]  $\rightarrow s'$   
 $S \rightarrow A \rightarrow$  [model]  $\rightarrow R$

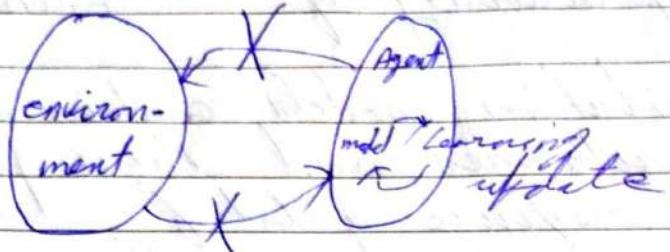
2. Q-learning update  $Q(S,A) \leftarrow Q(S,A) + \alpha [R + \gamma \max_a Q(s',a) - Q(S,A)]$

3. Greedy policy improvement

$$\pi(s) = \arg\max_a Q(S,a)$$

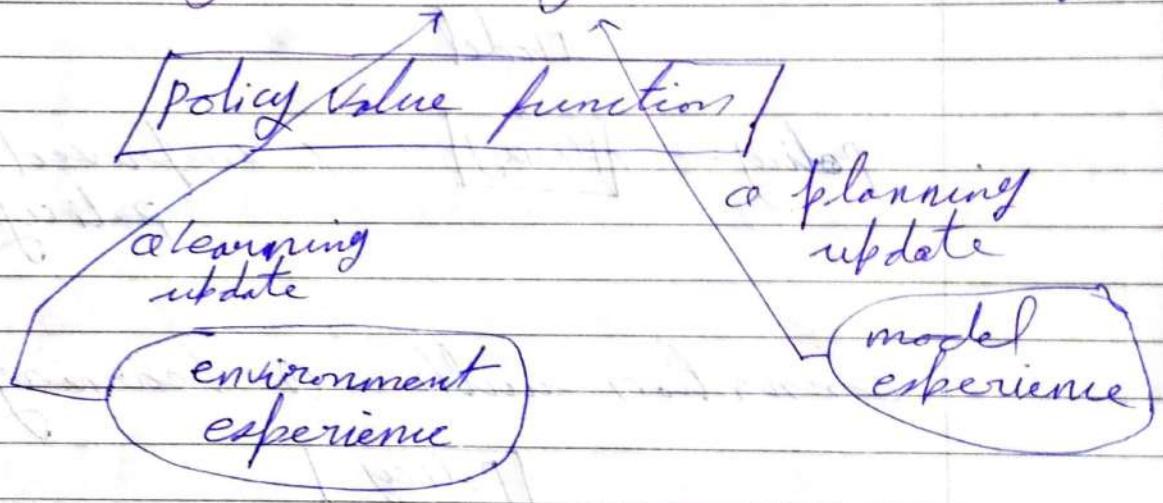
Date - / - / -

\* planning only uses imagined experience



-- The Dyna Architecture --

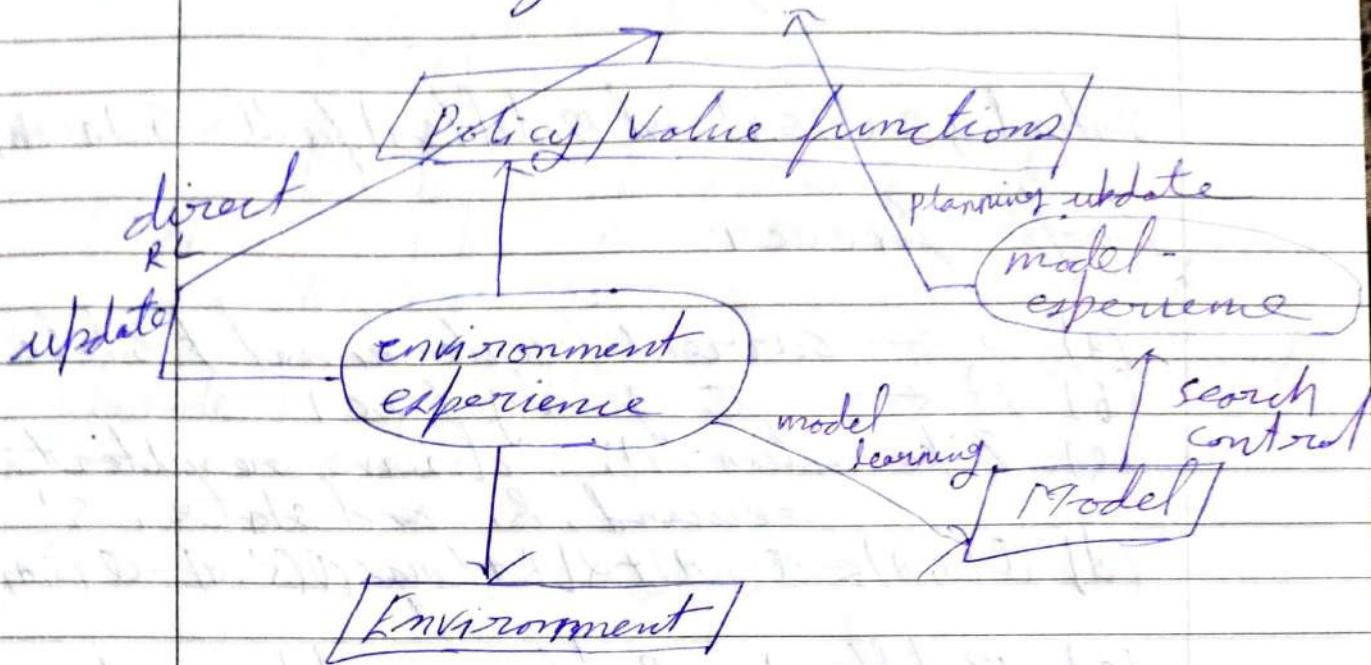
Combining Q-learning & Q-planning



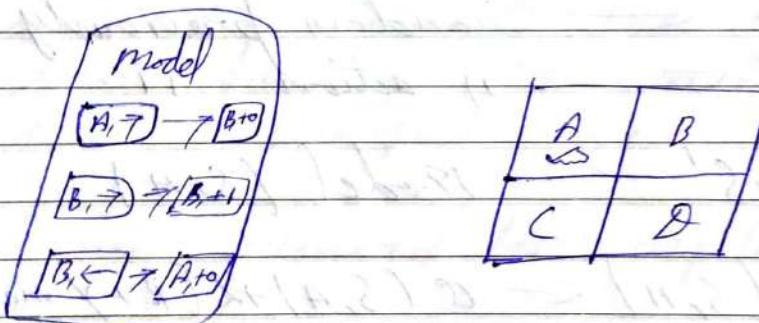
The environment experience can be used to learn the model. This model will be used to generate model experience. In addition, we want to control how the model generates this simulated experience, what states the agent will plan from. We call this process search control.

Date - 1 - 1 -

## The Dyna Architecture



## Learning a deterministic model



## Tabular Q-Learning

Initialize  $Q(s,a)$  and  $Model(s,a)$  for all  $s \in S, a \in A(s)$

Loop forever:

- (a)  $S \leftarrow$  current (nonterminal) state
- (b)  $A \leftarrow$   $\epsilon$  greedy ( $s, Q$ )
- (c) Take action  $A_i$ ; observe resultant reward,  $R$ , and state,  $S'$
- (d)  $Q(s, a) \leftarrow Q(s, a) + \alpha [R + \gamma \max_a Q(s', a) - Q(s, a)]$
- (e)  $Model(s, a) \leftarrow R, S'$  (assuming deterministic environment)
- (f) Loop repeat  $n$  times.

$s \leftarrow$  random previously observed state  
 $A \leftarrow$   $i$  action  $i$  taken in  $s$

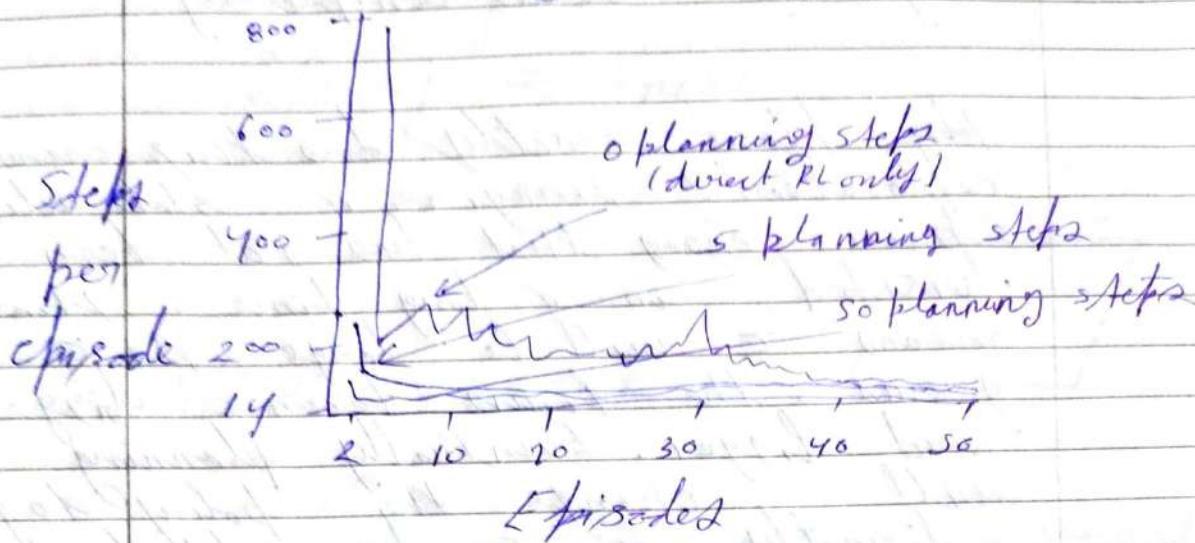
$R, S' \leftarrow Model(s, a)$

$Q(s, a) \leftarrow Q(s, a) + \alpha [R + \gamma \max_a Q(s', a) - Q(s, a)]$

In this algo, search control selects a previously visited state-action pair at random. It must be a state-action pair the agent has seen before, otherwise the model would not know what happens next.

Date - / - / -

More planning  $\rightarrow$  faster learning



Model inaccuracy: when the environment changes, the model will be incorrect. It will remain incorrect until the agent revisits that part of the environment that changed and updates the model.

Bonus reward for exploration

$$\text{New reward} = r + \kappa \sqrt{c}$$

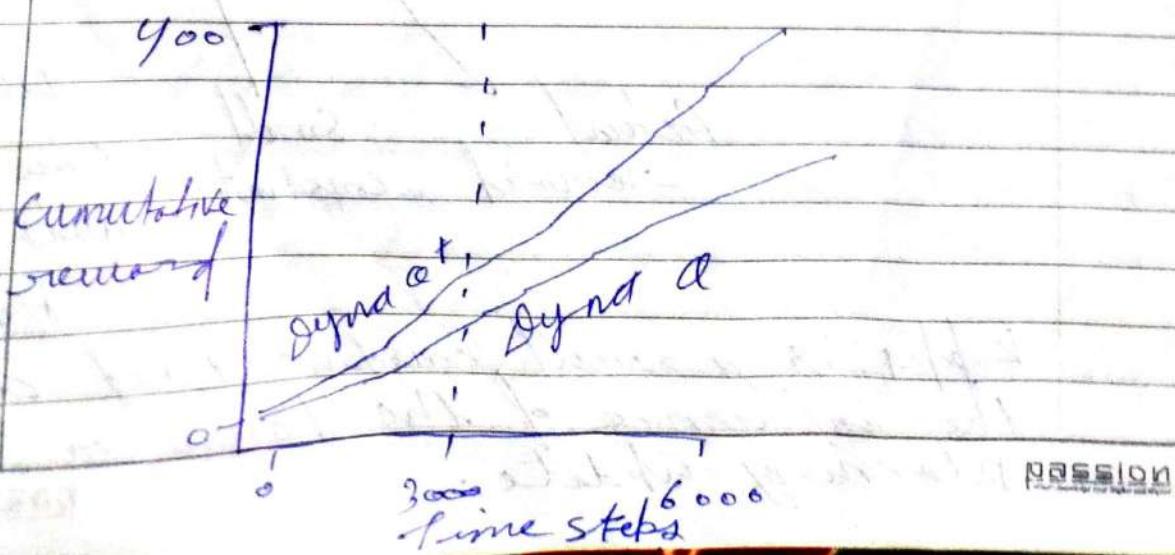
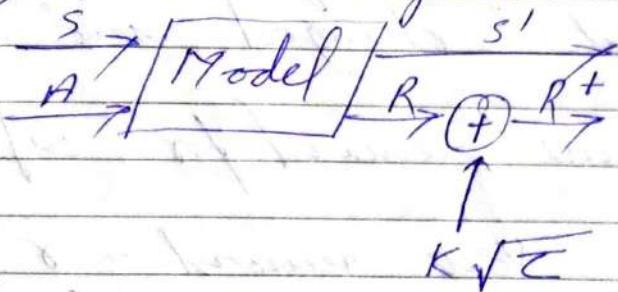
Actual reward  $\downarrow$  Small constant  $\uparrow$  time steps since transition was last stored

Kappa is a small constant that controls the influence of the bonus on the planning update.

If  $\kappa$  was zero, we would ignore bonus completely.

How exactly does  $\kappa$  encourage exploration? Imagine a state-action pair,  $(s, a)$ , that has not been visited in a long time. That means  $\tau$  will be large. As  $\tau$  grows, the bonus become bigger and bigger. Eventually planning will change the policy to go directly to  $s$  due to large bonus. When the agent finally visit state  $s$ , it might see a big reward or it might be disappointed.

The Dyna-Q+ algorithm



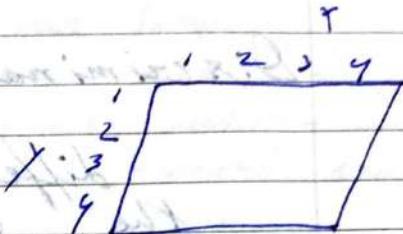
</>  
Course-3, Week 1 Date - 1-1-

parameterizing the value function

$$\hat{V}(s, w) \approx V_p(s)$$

weights

$$\hat{V}(s, w) = w_1 x + w_2 y$$



we only have to store the 2 weights

$$\hat{V}(s, w) = \sum w_i n_i(s)$$

This simply means that the value of each state is computed as the sum of the weights multiplied by some fixed attributes of the state called feature.

linear function approximation relies on having good features.

	High generalization
Low discrimination	High discrimination
High discrimination	Low generalization

Generalization  $\rightarrow$  we may not have to visit every state as much to get this values correct if we can learn its value from similar states.

Discrimination  $\rightarrow$  ability to make the values for two states different to distinguish b/w the values for these 2 states.

\* forming policy evaluation as supervised learning

$$\{ (S_1, \text{ctrl}), \dots, (S_1, \text{Retp}), (S_2, \text{ctrl}), \dots, (S_2, \text{Retp}), (S_2, \text{route}), (S_2, \text{car}), \dots \}$$

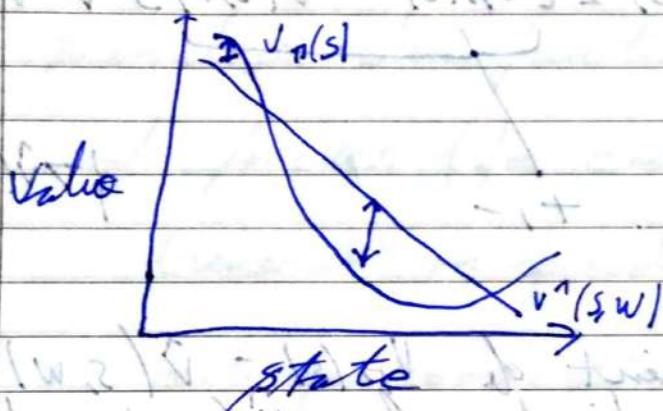
\* In reinforcement learning, an agent interacts with an environment and continually generates new data. This is often called the online setting. To distinguish it from the offline setting where the full dataset is available from the start and remains fixed throughout learning.

+ TD methods use bootstrapping.

meaning that our target now depends on our own estimate.

These estimation changes as learning progresses and so our target continually change. This is different than supervised learning where we have access to ground truth.

+ The Mean Squared Value Error



$$\text{Error} = \sum_{s \in S} [v_\pi(s) - v^*(s, w)]^2$$

$v_\pi(s)$  tells us how much we care about each state. A natural measure is the fraction of time each state is visited under the policy.

That the states that the policy spends more time in have a higher weight is the objective. We care

less about errors in the states the policy visits less frequently.

$$VE = \mathbb{E}_{s \sim S} [u(s) [V_p(s) - \hat{V}(s, w)]^2]$$

\* Gradient of the Mean Squared Value Error

$$\nabla_{w \in S} \sum_{s \in S} u(s) [V_p(s) - \hat{V}(s, w)]^2$$

$$= - \sum_{s \in S} u(s) 2 [V_p(s) - \hat{V}(s, w)] \nabla \hat{V}(s, w)$$

+,-

$$\nabla \hat{V}(s, w) = x(s)$$

\* The gradient of  $V_p(s) - \hat{V}(s, w)$  equals the gradient of  $\hat{V}(s, w)$  coz changing  $w$ - does not change  $V_p(s)$

\* If the difference is +ve, it means the true value is higher than our estimate, so we should change the weights in the direction that increases our estimate of the value if -ve, we should current error change the weights in opposite direction.

for complete,  $s_i$

$$W_2 = W_1 + \alpha [V_{\pi}(s_i) - \hat{V}(s_i, w_1)] \nabla \hat{V}(s_i, w_1)$$

Gradient Monte carlo algo :-

Input  $\pi$ : the policy  $\pi$  to be evaluated

Input  $\hat{V}$ : a differentiable function  $\hat{V}: S \times R^d \rightarrow R$

Algo parameters: step size  $\alpha > 0$

Initialize value function weights  $w \in R^d$

Loop forever (for each episode):

Generate an episode  $S_0, A_0, R_1, \dots, R_T, S_T$  using  $\pi$

loop for each step of episode  $t=0, \dots, T-1$

$$W \leftarrow W + \alpha [G_t - \hat{V}(S_t, w)] \nabla \hat{V}(S_t, w)$$

so  $\pi$   $\xleftarrow{\text{policy}}$   $\pi$ .  $\alpha = 2 \times 10^{-3}$

BT: twiddle Returns  $1, 1, 1, \dots, 1$  (twiddle)

Visited states  $500, 423, 482, \dots, 336$

each group  
100 states



T

Date - / - / -

The TD update for function approx -

$$w \leftarrow w + \alpha [U_t - V(s_t, w)] \nabla V(s_t, w)$$

$$U_t = R_{t+1} + \gamma V(s_{t+1}, w)$$

$U_t$  is biased  $\Rightarrow$  we may not converge to a local optimum.

because  $U_t$  (target) depends on our estimate of the value in the next state. This means our update could be biased because the estimate is our target may not be accurate.

loop for semi-gradient TD(0)

loop for each episode

choose  $A \sim \pi(\cdot | s)$

Take action  $A$ , observe  $R, s'$

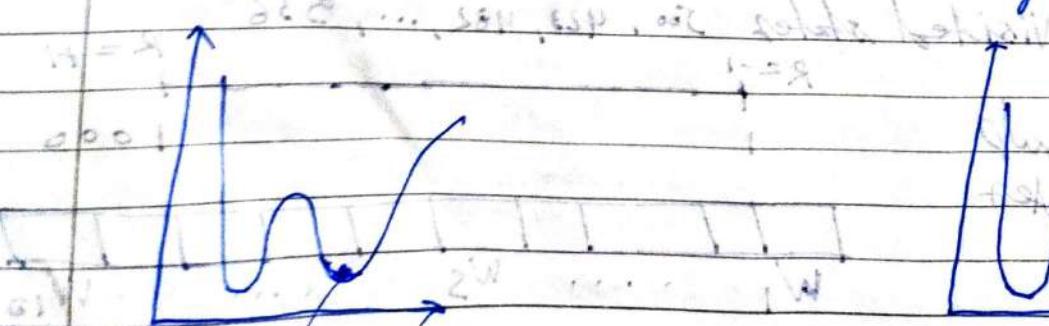
$$w \leftarrow w + \alpha [R + \gamma V(s', w) - V(s, w)] \nabla V(s, w)$$

$s \leftarrow s'$

until  $s$  is terminal

Gradient Monte Carlo

Semi-gradient TD

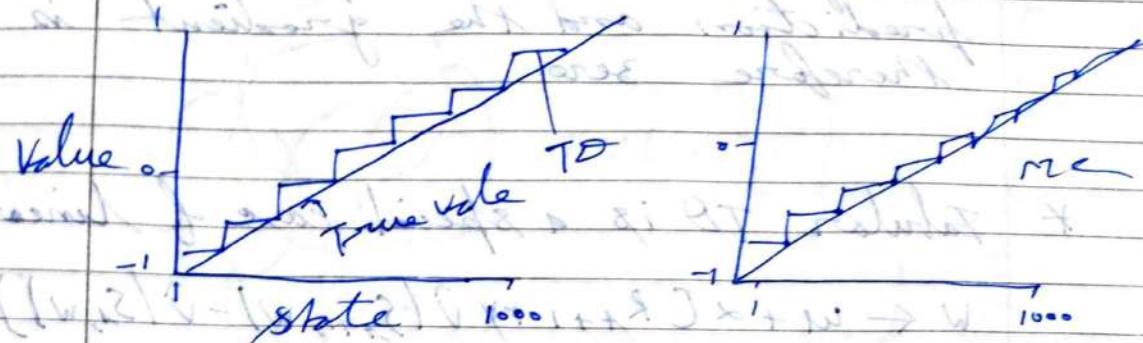


local minimum

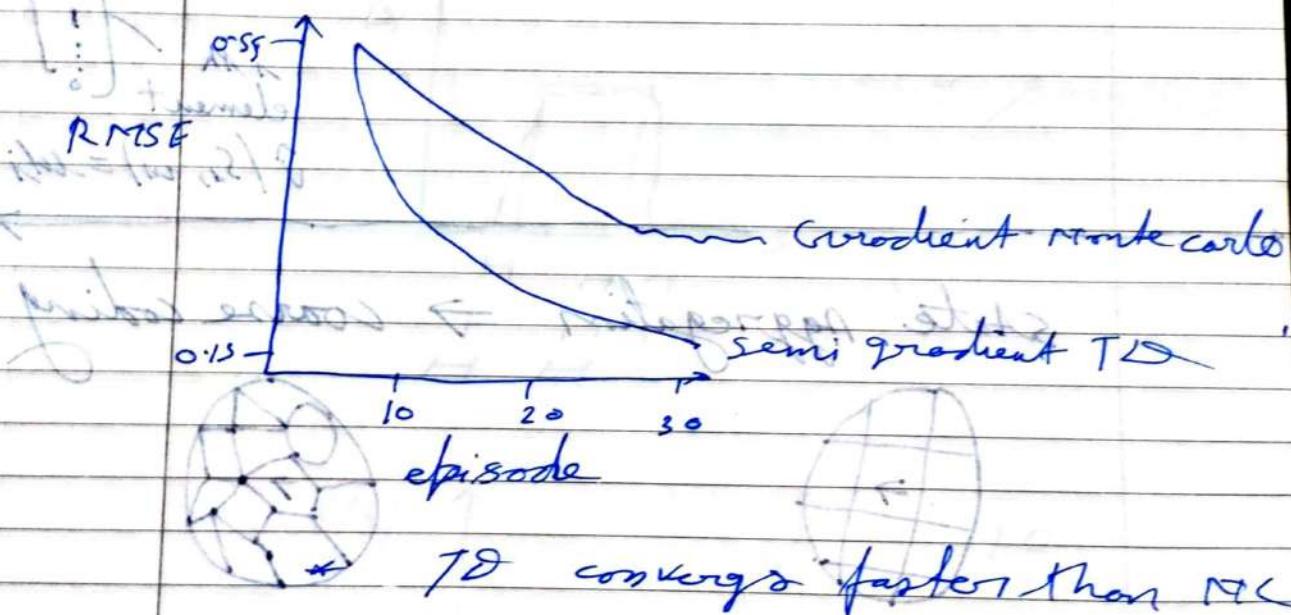
passion

Date - / - / -

## long run performance



$\text{TD} = \text{true value} + \alpha \text{TD}$  + MC better than TD



TD update with linear function app...

$$W \leftarrow w + \alpha S_t \nabla \hat{V}(S_t, w)$$
$$w \leftarrow w + \alpha S_t X(S_t)$$

$$\hat{V}(S_t, w) = w^T X(S_t)$$

$$\nabla \hat{V}(S_t, w) = X(S_t)$$

if the feature is zero, then that weight has no impact on the prediction and the gradient is therefore zero.

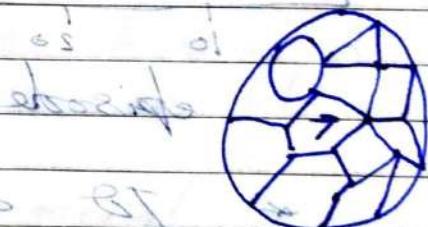
\* Tabular TD is a special case of linear TD

$$w \leftarrow w + \alpha [R_{t+1} + y \hat{v}(S_{t+1}, w) - \hat{v}(S_t, w)] x(S_t)$$

$$w_i \leftarrow w_i + \alpha [R_{t+1} + y \hat{v}(S_{t+1}, w) - \hat{v}(S_t, w)]_i y/S_i = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

*i*th element  
 $\hat{v}(S_t, w) = w_i$

~~state transition~~ ~~feature~~ ~~action~~ ~~reward~~ ~~observed~~ ~~state~~  $\rightarrow$  coarse coding



by allowing overlap, we obtain a many & flexible class of feature representations called coarse coding.

$$(w, b)v^T + b + w \rightarrow w$$

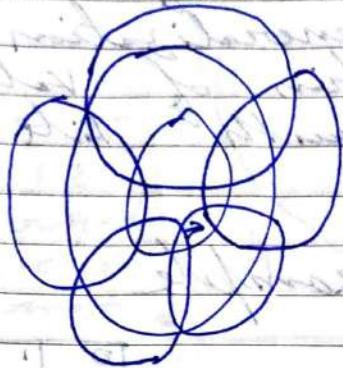
$$(b)v^T + b + w \rightarrow w$$

$$(b)v^T + w \approx (w, b)v^T$$

$$(b)v^T \approx (w, b)v^T$$

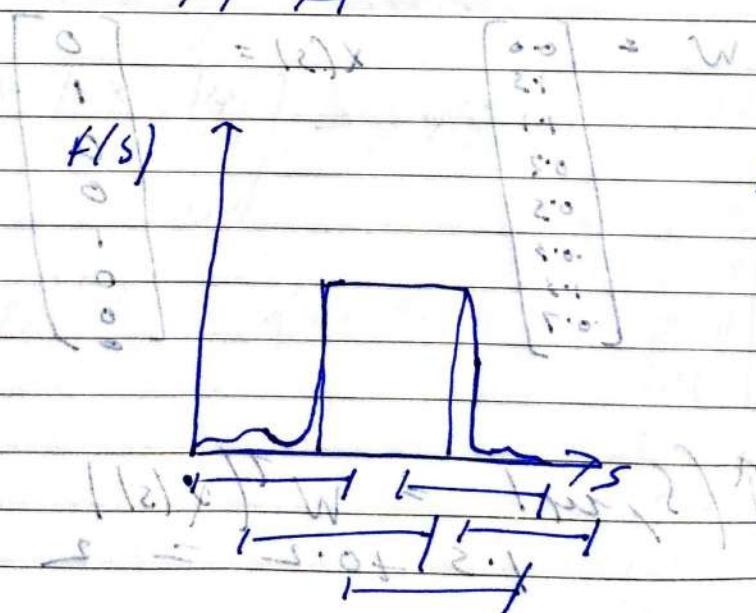
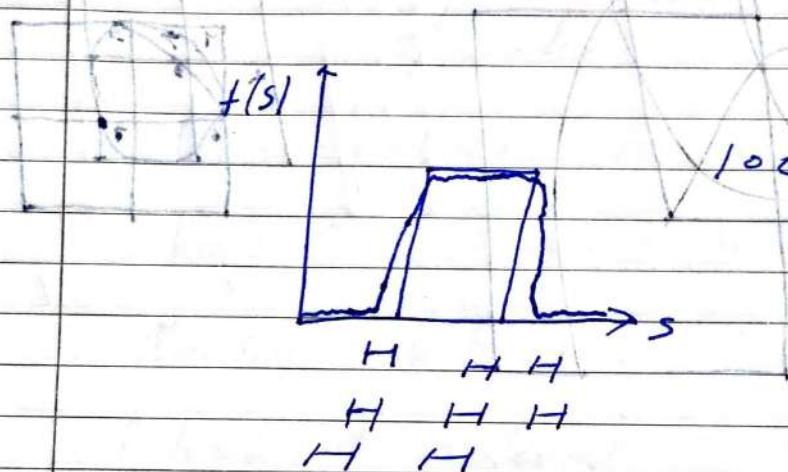
Date - / - / -

coarse coding



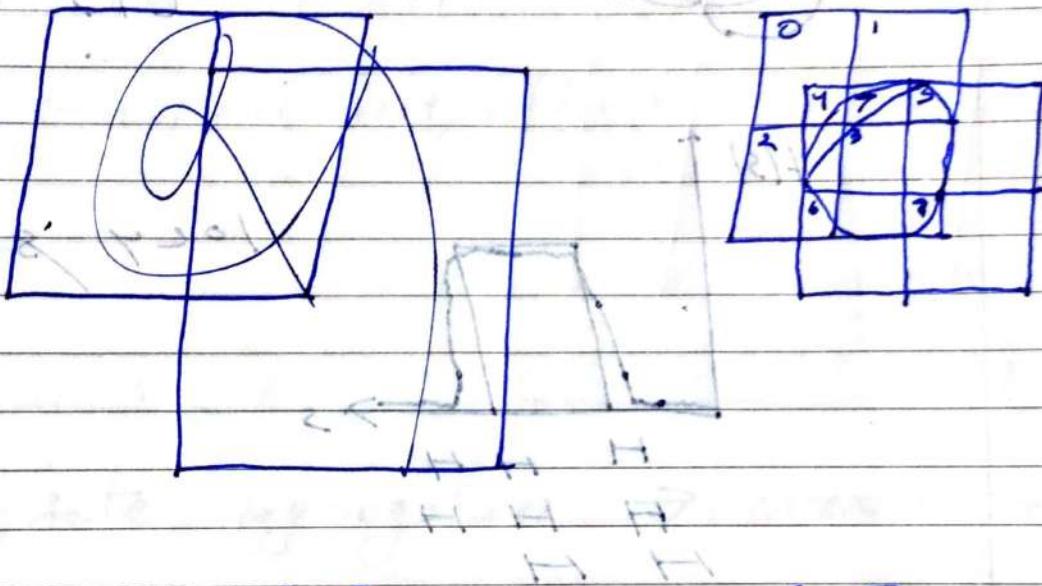
$$f(s) \rightarrow$$

0	0
0	1
0	0
0	0
1	0



\* The shape, size of tiles, and the number of tiling affects the generalization, and dissemination of value function approximation with tile coding.

A simple example

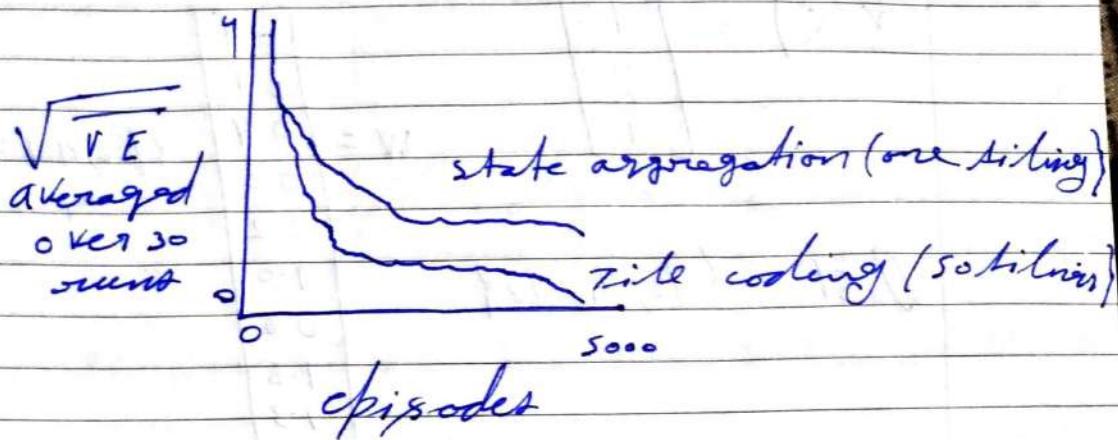


$$W = \begin{bmatrix} 0.0 \\ 1.5 \\ 1.1 \\ 0.2 \\ 0.5 \\ -0.3 \\ 1.3 \\ -0.7 \end{bmatrix} \quad x(s) = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

$$V(s, a) = W^T(x(s)) \\ 1.5 + 0.2 = 2$$

Date - 1 - 1 -

## Tile coding vs state aggregation



state-values to action-values

$$V_\pi(s) \approx \hat{V}(s, \pi) = W^T x(s)$$

$$q_\pi(s, a) \approx \hat{q}(s, a, \pi) = W^T x(s, a)$$

Representing actions

$$x(s) = \begin{bmatrix} x_0(s) \\ x_1(s) \\ x_2(s) \\ x_3(s) \end{bmatrix} \rightarrow x(s, a) = \begin{bmatrix} u_0(s) \\ n_1(s) \\ n_2(s) \\ n_3(s) \end{bmatrix} \left\{ \begin{array}{l} a_0 \\ a_1 \\ a_2 \end{array} \right\} \begin{bmatrix} u_0(s) \\ n_1(s) \\ n_2(s) \\ n_3(s) \end{bmatrix}$$

$$\mathcal{A}(s) = \{a_0, a_1, a_2\} \quad x(s, a_0) = \begin{bmatrix} u_0(s) \\ n_1(s) \\ n_2(s) \\ n_3(s) \end{bmatrix} \left\{ \begin{array}{l} a_0 \\ a_1 \\ a_2 \end{array} \right\} \begin{bmatrix} u_0(s) \\ n_1(s) \\ n_2(s) \\ n_3(s) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$x(s, a_1) = \begin{bmatrix} u_0(s) \\ n_1(s) \\ n_2(s) \\ n_3(s) \end{bmatrix} \left\{ \begin{array}{l} a_0 \\ a_1 \\ a_2 \end{array} \right\} \begin{bmatrix} u_0(s) \\ n_1(s) \\ n_2(s) \\ n_3(s) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$x(s, a_2) = \begin{bmatrix} u_0(s) \\ n_1(s) \\ n_2(s) \\ n_3(s) \end{bmatrix} \left\{ \begin{array}{l} a_0 \\ a_1 \\ a_2 \end{array} \right\} \begin{bmatrix} u_0(s) \\ n_1(s) \\ n_2(s) \\ n_3(s) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$x(s, a_0) = \begin{bmatrix} u_0(s) \\ n_1(s) \\ n_2(s) \\ n_3(s) \end{bmatrix} \left\{ \begin{array}{l} a_0 \\ a_1 \\ a_2 \end{array} \right\} \begin{bmatrix} u_0(s) \\ n_1(s) \\ n_2(s) \\ n_3(s) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$x(s, a_1) = \begin{bmatrix} u_0(s) \\ n_1(s) \\ n_2(s) \\ n_3(s) \end{bmatrix} \left\{ \begin{array}{l} a_0 \\ a_1 \\ a_2 \end{array} \right\} \begin{bmatrix} u_0(s) \\ n_1(s) \\ n_2(s) \\ n_3(s) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$x(s, a_2) = \begin{bmatrix} u_0(s) \\ n_1(s) \\ n_2(s) \\ n_3(s) \end{bmatrix} \left\{ \begin{array}{l} a_0 \\ a_1 \\ a_2 \end{array} \right\} \begin{bmatrix} u_0(s) \\ n_1(s) \\ n_2(s) \\ n_3(s) \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Date - / - / -

## computing action-values

$$x(s_0) = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

$$w = \begin{bmatrix} 0.7 \\ 0.1 \\ 0.4 \\ 0.3 \\ 1.2 \\ 1.0 \\ 0.6 \\ 1.8 \\ 1.3 \\ 1.1 \\ 0.9 \\ 0.7 \end{bmatrix}$$

$$x(s_0)a_0 = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \\ 0.7 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$A(S) = \{a_0, a_1, a_2\}$$

$$\hat{q}(s_0, a_0, w) = 0.7 + 0.3 = 1$$

$$\hat{q}(s_0, a_1, w) = 2.2 + 1.8 = 4$$

$$\hat{q}(s_0, a_2, w) = 1.3 + 1.7 = 3$$

Date - 1 - 1 -

+ Episodic Sarsa with function approx.

Episodic semi-gradient Sarsa for estimating

$$\hat{Q} = Q_+$$

Input: a differentiable action value function parameterization  $\hat{Q}: S \times A \times R^d \rightarrow \mathbb{R}$

Algo parameters: step size  $\alpha > 0$ , small  $\epsilon > 0$

Initialize value function weights  $w \in \mathbb{R}^d$   
arbitrarily (e.g.,  $w=0$ )

Loop for each episode:

$s, a \leftarrow$  initial state and action of episode  
| e.g.:  $\epsilon$ -greedy

Loop for each step of episode:

Take action  $A$ , observe  $R, s'$   
if  $s'$  is terminal:

$$w \leftarrow w + \alpha [R + \gamma \hat{Q}(s', a', w) - \hat{Q}(s, a, w)] \nabla \hat{Q}(s, a, w)$$

Go to next episode

choose  $A'$  as a function of  $\hat{Q}(s', \cdot, w)$   
| e.g.,  $\epsilon$ -greedy

$$w \leftarrow w + \alpha [R + \gamma \hat{Q}(s', a', w) - \hat{Q}(s, a, w)] \nabla \hat{Q}(s, a, w)$$

$$s \leftarrow s'$$

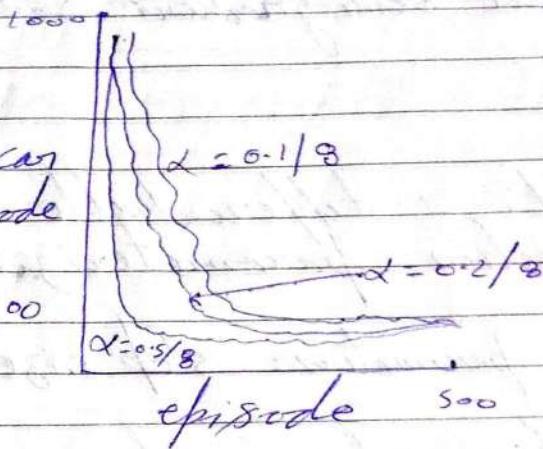
$$A \leftarrow A'$$

Date - / - / -

## Learning curves

Mountain car  
steps per episode

Averaged over 100  
runs



### Expected SARSA with function Approximation

SARSA :

$$w_t \leftarrow w_t + \alpha [R_{t+1} + \gamma \hat{q}(S_{t+1}, a_{t+1}) - \hat{q}(S_t, A_t, w_t)] \nabla \hat{q}(S_t, A_t, w_t)$$

Expected SARSA :

$$\begin{aligned} w_t &\leftarrow w_t + \alpha [R_{t+1} + \gamma \sum_{a'} \pi(a'|s_{t+1}) \hat{q}(s_{t+1}, a', w_t) \\ &\quad - \hat{q}(s_t, A_t, w_t)] \nabla \hat{q}(s_t, A_t, w_t) \end{aligned}$$

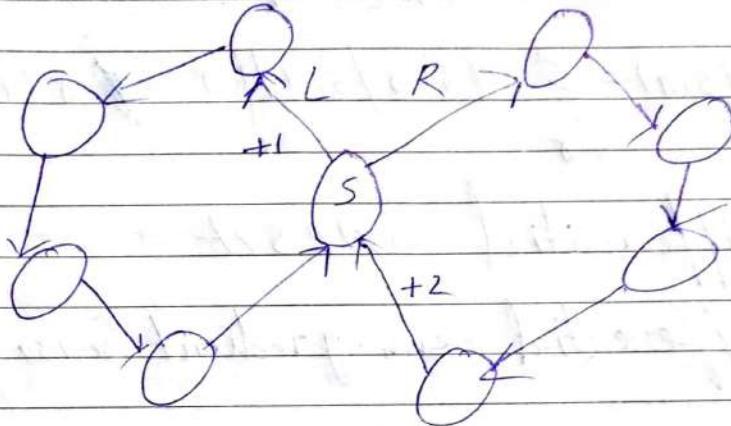
O Learning

$$w_t \leftarrow w_t + \alpha [R_{t+1} + \gamma \max_{a'} \hat{q}(s_{t+1}, a', w_t) - \hat{q}(s_t, A_t, w_t)] \nabla \hat{q}(s_t, A_t, w_t)$$

Date - / - / -

## The Average Reward objective

$$\delta(\pi) = \sum_{s} \sum_{a} \pi(a|s) \sum_{s',r} p(s',r|s,a) r$$



$$\delta[\pi_L] = 1/s = 0.2$$

$$\delta[\pi_R] = 2/s = 0.4$$

Returns for Average Reward

$$G_t = R_{t+1} - \delta(\pi) + R_{t+2} - \delta(\pi) + R_{t+3} - \delta(\pi) + \dots$$

$$\begin{aligned} \pi(\pi_L) &\rightarrow G_t = 1 - 0.2 + \\ &\quad 0 - 0.2 + \\ &\quad \vdots H \geq 0 \\ &\quad = 0.4 \end{aligned}$$

$$\begin{aligned} G_t &= 0 - 0.2 + \\ &\quad 2 - 0.2 + \\ &\quad \vdots \\ &\quad = 0.4 \end{aligned}$$

$$\begin{aligned} \pi(\pi_R) &\rightarrow G_t = 1 - 0.4 + \\ &\quad 0 - 0.4 + \\ &\quad 0 - 0.4 + \\ &\quad \vdots \\ &\quad = 0.8 \end{aligned}$$

$$\begin{aligned} G_t &= 0 - 0.4 + \\ &\quad 2 - 0.4 + \\ &\quad \vdots \\ &\quad = 0.8 \end{aligned}$$

Date - 1-1

## Value functions for Average Reward

$$G_t = R_{t+1} - \gamma/\pi + R_{t+2} - \gamma^2/\pi + \dots$$

$$q_\pi(s, a) = E_\pi[G_t | S_t = s, A_t = a]$$

$$Q_\pi(s, a) = \sum_{s', a'} P(s', a'|s, a) [r - \gamma \pi(a'|s) + Q_\pi(s', a')]$$

### \* Differential SARSA

Differential semi-gradient SARSA for estimating

$$\hat{Q} \leftarrow \hat{Q} + \eta_x$$

Input:  $\hat{Q}: S \times A \times R^d \rightarrow R$

Algo params: step sizes  $\alpha, \beta > 0$

Initialize value function, weights  $w \in R^d$  and  
average reward estimate  $\bar{R} \in R$

Initialize states  $s$ , and actions  $A$

Loop for each step:

Take action  $A$ , observe  $R, s'$

Choose  $A'$  as a function of  $\hat{Q}(s, a, w)$

$$\delta \leftarrow R - \bar{R} + \hat{Q}(s, A, w) - \hat{Q}(s, A', w)$$

$$\bar{R} \leftarrow \bar{R} + \beta \delta$$

$$w \leftarrow w + \alpha \delta \nabla \hat{Q}(s, A, w)$$

$$s \leftarrow s'$$

$$A \leftarrow A'$$

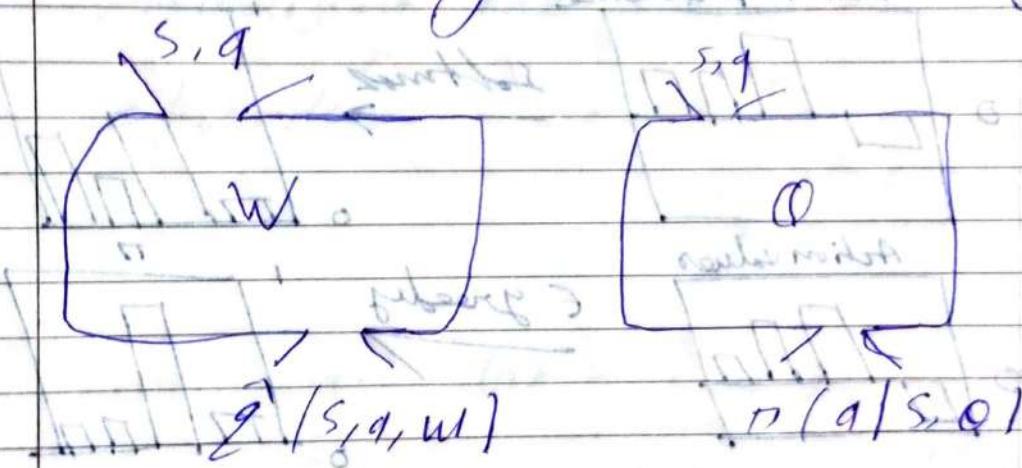
Date - 1-1-

\* An example of a policy Without action-values

Velocity |  $\pi(\text{Right}|s) = 1$   
 $\pi(\text{Left}|s) = 0$

position

parametrizing policies Directly



The Softmax policy parameterization

$$\pi(a|s, \theta) = \frac{e^{h(s, a, \theta)}}{\sum_i e^{h(s, i, \theta)}}$$

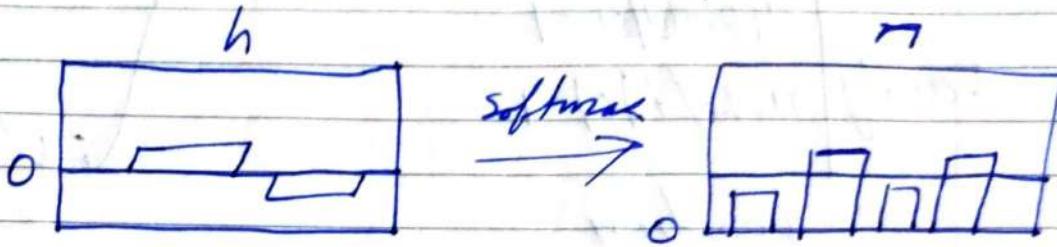
action preference

$h(s, a, \theta)$  =  $\theta^\top \phi(s, a)$

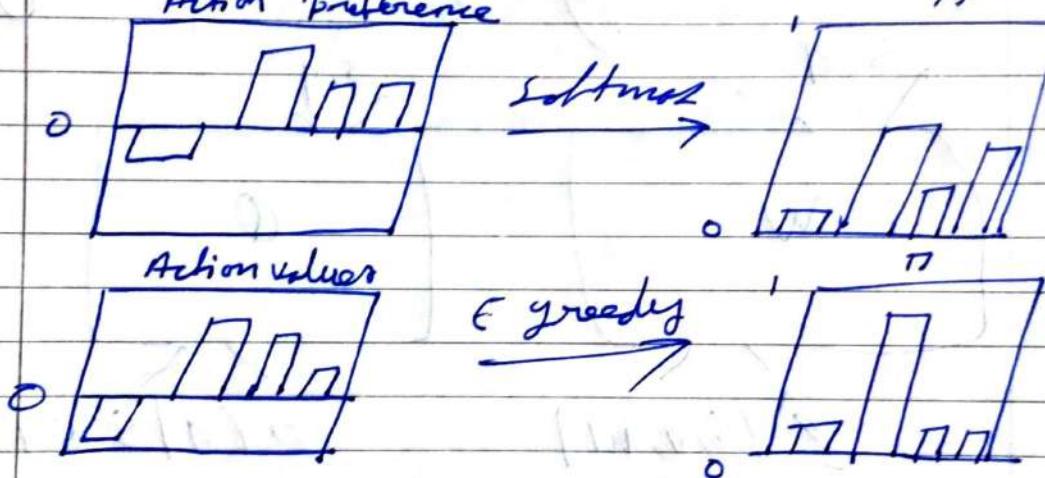
of softmax scores

## ~~Transpose Convolution~~

- \* Examples of softmax policies



- \* Action preferences are not action values



- \* Advantages of policy parameterization

Epsilon step chooses a random action to ensure continual exploration. However,  $\epsilon$  probabilities puts a cap on how good the resulting policy can be.

we could f. coarse switch to greedy policy when

Date - / - / -

exploration is done.

We can avoid this issue with parameterized policies. The policy can start off stochastic to guarantee expiration. Then as learning progresses, the policy can naturally converge towards a deterministic greedy policy.

\* Formalizing the goal as an objective

$$\text{Episodic: } G_t = T$$

$$\sum_{t=0}^T R_t$$

$$\text{Continuing: } G_t = \sum_{t=0}^{\infty} \gamma^t R_t = \sum_{t=0}^{\infty} R_t \gamma^{t-1}$$

The Average Reward objective

$$\gamma(\pi) = \mathbb{E}_{\pi}[s] \mathbb{E}_\pi \pi(a|s, 0) \sum_s \sum_a \mathbb{P}(s', r|s, a) \gamma$$
$$E_\pi[\lambda_t] \quad \downarrow \quad E_\pi[R_t | S_t = s]$$
$$E_\pi[R_t | S_t = s, A_t = a]$$

Date - / - / -

$$\nabla \delta(\pi) = \nabla \mathbb{E}_{\pi(s)} [\sum_a \pi(a|s, \theta) \mathbb{E}_{p(s, a|s, \theta)} r]$$

$\uparrow s$        $\uparrow a$        $\uparrow s, \theta$

depends on  $\theta$

$$\nabla \delta(\pi) = \sum_a \pi(a) \mathbb{E} [r(s, a) \nabla \pi(a|s, \theta)]$$

\* Unbiasedness of the stochastic samples

$$\nabla \delta(\pi) = \mathbb{E}_\pi \left[ \sum_a \mathbb{E} [r(s, a) \nabla \pi(a|s, \theta)] \right]$$

$s_0, A_0, R_1, S_1, A_1, \dots, S_t, A_t, R_t, \dots$

Getting stochastic samples with one action

$$= \sum_a \nabla \pi(a|s, \theta) \mathbb{E}_{\pi} [r(s, a)]$$

$$= \sum_a \pi(a|s, \theta) \frac{\nabla \pi(a|s, \theta)}{\pi(a|s, \theta)} \mathbb{E}_{\pi} [r(s, a)]$$

$$= \mathbb{E}_\pi \left[ \frac{\nabla \pi(a|s, \theta)}{\pi(a|s, \theta)} \mathbb{E}_{\pi} [r(s, a)] \right]$$

Stochastic Gradient Ascent for policy  
Date - 1 - parameter

$$\theta_{t+1} = \theta_t + \alpha \frac{\nabla \ln \pi(A_t | S_t, \theta_t) Q_{\pi}(S_t, A_t)}{\nabla \ln \pi(A_t | S_t, \theta_t)}$$

$$\theta_{t+1} = \theta_t + \alpha \nabla \ln \pi(A_t | S_t, \theta_t) Q_{\pi}(S_t, A_t)$$

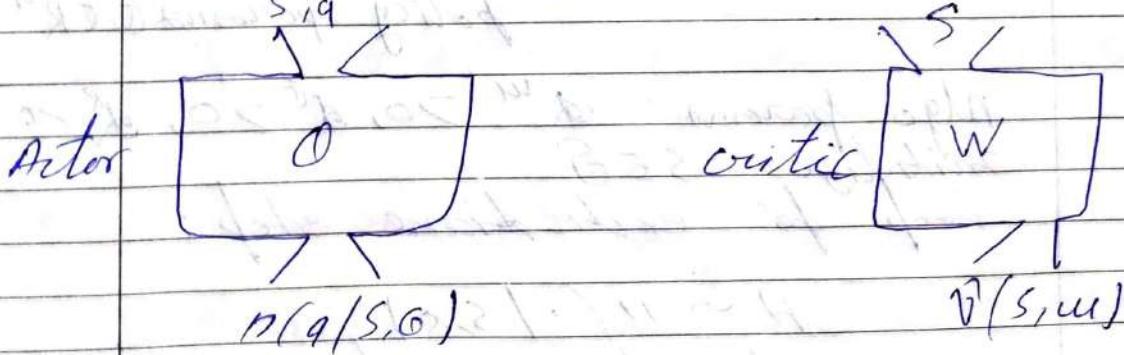
$$\nabla \log(\pi(a|s)) = \frac{\nabla f(s)}{f(s)}$$

~~Approximating~~

\* Actor - critic

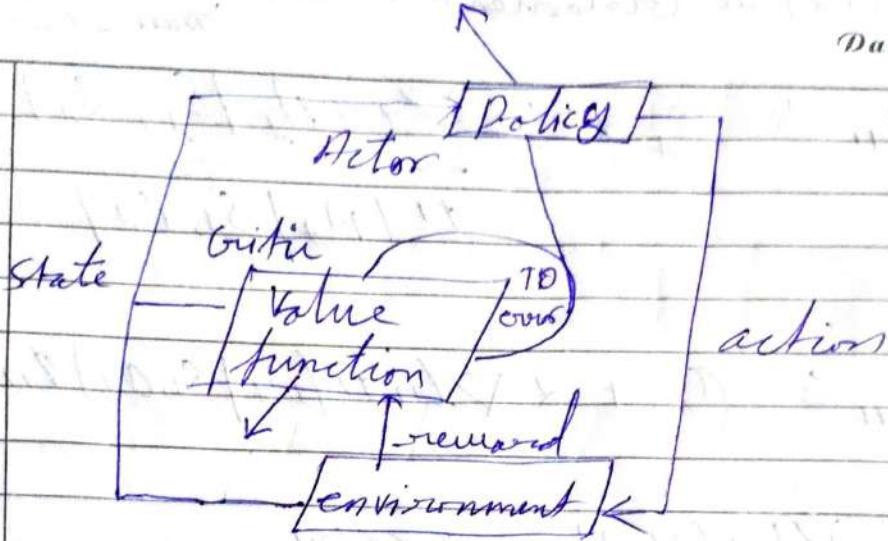
approximating the action value in policy update

$$\theta_{t+1} = \theta_t + \alpha \nabla \log \pi(A_t | S_t, \theta_t) [R_{t+1} - \bar{V}(S_{t+1}, w)]$$



\* The actor is continually changing the policy to exceed the critics expectation, and the critic is constantly updating its value function to evaluate the actor changing policy.

Date - / - / -



Actor - Critic (continuing), for estimating  $\pi(a|s, o)$

Input: a differentiable policy parametrization  $\pi(a|s, o)$

Input: a "state-value function parametrization  $V(s, u)$

Initialize  $\bar{R} \in \mathbb{R}$  to 0

Initialize state-value weights  $w \in \mathbb{R}^d$  &  
policy params  $\theta \in \mathbb{R}^d$

Algo params:  $\alpha^w > 0, \alpha^0 > 0, \alpha^R > 0$

Initialize  $s \in S$

Loop for each time step:

$$A \sim \pi(\cdot | s, o)$$

Take action  $A$ , observe  $S', R, R'$

$$S' \leftarrow R - \bar{R} + \hat{V}(s, u) - V(s, u)$$

$$\bar{R} \leftarrow \bar{R} + \alpha^R S'$$

$$u \leftarrow u + \alpha^w S' \nabla V(S, u)$$

$$\theta \leftarrow \theta + \alpha^0 S' \nabla \log \pi(A | s, o)$$

$$s \leftarrow s'$$

Date - / - / -

\* Actor-Critic with Softmax policies

$$\theta \leftarrow \theta + \alpha^{\theta} S \nabla \log \pi(A|S, \theta)$$

$$\pi(a|S, \theta) = \frac{e^{h(S, a, \theta)}}{\sum_{b \in A} e^{h(S, b, \theta)}}$$

$$\hat{V}(S, w) = w^T x(S)$$

$$h(S, a, \theta) = \theta^T x_h(S, a)$$

$$w \leftarrow w + \alpha^w S \nabla \hat{V}(S, w)$$

$$\nabla \hat{V}(S, w) = x(S)$$

$$\theta \leftarrow \theta + \alpha^{\theta} S \nabla \log \pi(A|S, \theta)$$

$$w \leftarrow w + \alpha^w S x(S)$$

$$\nabla \log \pi(a|S, \theta) = x_h(S, a) - \sum_b \pi(b|S, \theta) x_h(S, b)$$

\* Gaussian policy

$$\pi(a|S, \theta) = \frac{1}{\sigma(S, \theta) \sqrt{2\pi}} \exp\left(-\frac{(a - \mu(S, \theta))^2}{2\sigma^2(S, \theta)}\right)$$

$$\mu(S, \theta) = \theta_a^T x(S)$$

$$\sigma(S, \theta) = \exp(\theta_\sigma^T x(S))$$

$$\theta = \begin{bmatrix} \theta_u \\ \theta_c \end{bmatrix}$$

\* Gradient of the log of the Gaussian Policy

$$\nabla \log \pi(a | s, \theta_u) = \frac{1}{\sigma(s, \theta)^2} (a - u(s, \theta)) x(s)$$

$$= \left| \frac{(a - u(s, \theta))^2}{\sigma(s, \theta)^2} - 1 \right| x(s)$$

\* Advantages of Continuous Actions

- (1) It might not be straight forward to choose a proper discrete set of actions
- (2) Continuous action allows us to generalize over actions

For example, the gaussian policies smoothly assigns probability density to nearby actions. Imagine an action is found to be good and the update increases density for that action.

The density for a nearby action will also increase with the agent generalizing that those actions are likely to be good.

