

RL Final

Epsilon soft policy, On policy / Off policy

On policy - behavioral policy \rightarrow target policy

(take action)

(learn)

.. Epsilon soft, Epsilon greedy

non-optimal

Epsilon soft is good in case of behavioral policy because it gives us exploration

but not in case of target because epsilon soft stops us from taking greedy action everytime

Epsilon soft is not bad because exploration is necessary too but it is also not optimal

Off policy

behavioral policy \neq target policy



(M|T|W|T|F|S)

Date

Policy

probability of taking action

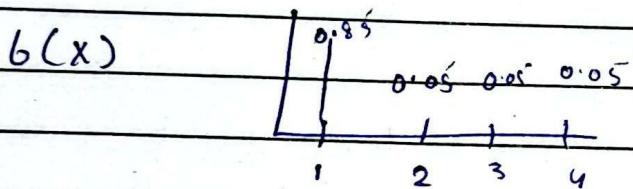
behavioral policy

Target policy

= probability distribution

e.g.

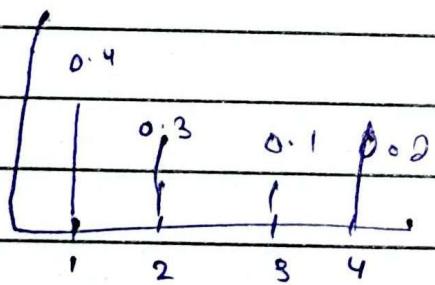
$X \rightarrow$ random variable



$$E(X) = \sum_{x \in X} x \cdot b(x)$$

$$= 1(0.85) + 2(0.05) + 3(0.05) + 4(0.05)$$

$\pi(x)$



$$= 1(0.4) + 2(0.3) + 3(0.1) + 4(0.2)$$

If probability is not given

we

repeat experiment for a large
no. of time and calculate the average.

Date _____

[M T W T F S]



Now we want to determine the estimated value of v under π by using samples from B :

Now we know the choice of b and π is different.

We modify the values of B such a way that we get expected value of x .

Important Sampling

(sample generate from behavior policy but we have to modify their return so it sets equal to the expected value of target policy)

$$E_{\pi}(x) = \mathbb{E}_{x \in X} u \cdot \pi(x)$$

$$= \mathbb{E}_{x \in X} \frac{u \cdot b(x)}{b(x)} \pi(x)$$

$$= \left(\mathbb{E}_{x \in X} \frac{u \cdot b(x)}{b(x)} \right) \left(\frac{\pi(x)}{b(x)} \right) \text{ outcome } k$$

Expected value Importance sampling
How much importance is given to k

$$= E_b(x) \cdot \rho(x)$$

$$= \frac{1}{n} \sum_{i=1}^n x_i \cdot \rho(x)$$

average under

Date _____

MTWTF



for above example

$$1 \times \frac{0.4}{0.85} + 1 \times \frac{0.4}{0.85} + 1 \times \frac{0.4}{0.85} +$$

$$2 \times \frac{0.3}{0.05} + 3 \times \frac{0.1}{0.05}$$

no of samples

How to Find $p(n)$ in Monte Carlo $\mathcal{Q} \rightarrow \Pr(\text{trajectory under } \pi) \rightarrow \text{target}$ $\Pr(\text{trajectory under } \theta) \rightarrow \text{below}$

Importance sampling estimates the value functions for a policy π with samples collected previously from an older policy θ'



Date _____

(M|T|W|T|F|S)

How to find important sampling in Monte Carlo

$$\rho = \frac{\Pr(\text{trajectory under } \pi)}{\Pr(\text{trajectory under } b)} \rightarrow \text{target}$$

$$\Pr(\text{trajectory under } b) \rightarrow \text{behavior}$$

$$\begin{aligned} \Pr(\text{trajectory under } \pi) &= \Pr(S_0, A_0, S_1, A_1, S_2, A_2, S_3, \dots, A_T, S_T) \\ &= \pi(A_0 | S_0) p(S_1 | S_0, A_0) \end{aligned}$$

$$= \Pr(A_t, S_{t+1}, A_{t+1}, S_{t+2}, A_{t+2}, \dots, A_T, S_T | S_0, A_0)$$

$$= \pi(A_t | S_t) p(S_{t+1} | S_t, A_t)$$

$$= \pi(A_{t+1} | S_{t+1}) p(S_{t+2} | S_{t+1}, A_{t+1})$$

terminal state

$$\Pr(\text{trajectory under } \pi) = \prod_{k=0}^{T-1} \pi(A_k | S_k) p(S_{k+1} | S_k, A_k)$$

$$\Pr(\text{trajectory under } b) = \prod_{k=0}^{T-1} \pi(A_k | S_k) p(S_{k+1} | S_k, A_k)$$

S_0

$$\rho_{t:T-1} = \frac{\pi(A_k | S_k)}{p(A_k | S_k)}$$

18 / 19

Date 20, 21, 22, 23

(M/T/W/T/F/S)



$$C_{t-1} = R_t + C_t$$

$$C_{t-1} = \frac{\pi(A_{t-1} | S_{t-1})}{b(A_t | S_t)}$$

$$G_{T-1} = R_{t+1}$$

You have to
multiply your
samples with
a value
so that
you modify
your behavior
policy towards
target policy

$$G_{T-1} * R_{T-1}$$

$$G_{T-2} = R_{T-1} + R_{T-2}$$

$$= G_{T-1} + R_{T-2}$$

$$G_{T-2} * R_{T-2}$$

$W \rightarrow$ importance sampling ratio

On-policy \rightarrow Off-policy

Multiply sample by an importance sampling ratio

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



M T W T F S

Date _____

Temporal Difference (Sample based)

$V(S_t)$ - Average (Returns)

= Average (Returns[0], Returns[1], ...,)

$$V(S_t) = \frac{1}{n} \sum_{i=1}^{n-1} G_i$$

MATB

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

$$= \frac{1}{n} \left[\sum_{i=1}^{n-1} R_i + R_n \right]$$

$$= \frac{1}{n} \left[\frac{n-1}{n-1} \sum_{i=1}^{n-1} R_i + R_n \right]$$

$$= \frac{1}{n} \left[(n-1) Q_n + R_n \right]$$

$$= \frac{1}{n} \left[n Q_n - Q_n + R_n \right]$$

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

New estimate = old + step size [Target-old value]

Incremental
update in
Monte Carlo

$$V(S_t) = \frac{1}{n} \sum_{i=1}^n G_i = V(S_t) + \alpha [G_t - V(S_t)]$$

Problem: We can't find $V(s)$ until the episode is finished

(30) ————— 5 ————— (35) ————— (15) ————— (10) ————— (1)

$v(s_i)$

Date _____

[M T W T F S]



$$v(s_i) = v(s_i) + \alpha [G_{t+1} - v(s_i)]$$

In bellman equation we can estimate $v(s)$ by from the next $(v(s))$

DP (Bootstrapping) based on some values you can estimate future values.

Temporal Difference — { Monte Carlo — Sampling
PP — Bootstrapping

Bellman

$$V_t(s_t) = \mathbb{E}(G_t | S_{t+1} = s)$$

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

$$V_t(s_t) = \mathbb{E}(R_{t+1} + \gamma G_{t+1} | S_{t+1} = s)$$

$$V_t(s_t) = \mathbb{E}(R_{t+1} + \gamma V(s_{t+1}, S_{t+2}, s))$$

$$G_t \approx R_{t+1} + \gamma V(s_{t+1})$$

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



Date _____

M T W T F S

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

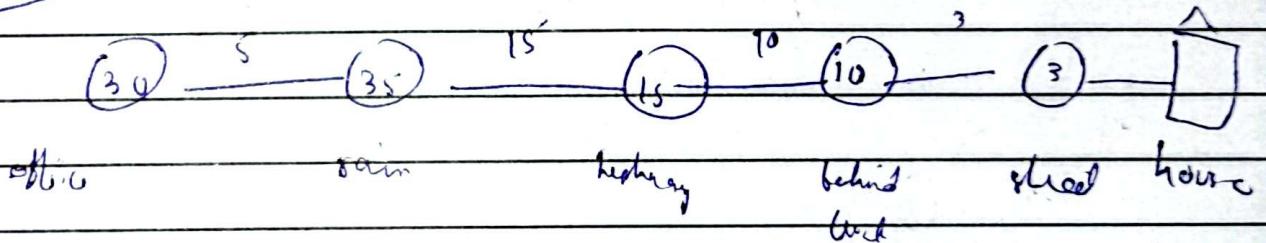
$$G_t \approx R_{t+1} + \gamma V(S_{t+1})$$

$$G_t \approx R_{t+1} + \gamma V(S_{t+1})$$

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

Temporal Difference

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



$$V(S_1) = 30 + 1 [5 + 1(35) - 30] \quad 30 + 5 = 35$$

$$V(S_2) = 35 + 1 [15 + 1(15) - 35] \quad 35 + 5 = 30$$

Tabular TD (0)
Tabular looking one step ahead.

number of states is finite



Date _____

M T W T F S

Intro to Q planning and DynaQ

Multi-armed bandits

Learning Estimate $\rightarrow Q_t(A_i) = \frac{1}{t-1} \sum_{i=1}^t R_i$

MDP:

Model

$$\text{based on } V_\pi(s) = \sum_a \pi(a|s) \sum_s \sum_r p(s', r | s, a) [r + \gamma V_\pi(s')]$$

Planning

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

Controls:

policy iteration
value iteration

Model free learning

Monte Carlo - Prediction

- control - on-policy
- off policy

TD - prediction

control \rightarrow SARSA

Date _____

(M T W T F S)



$$Q(s, a) = Q(s, a_t) + \alpha [R_t + \gamma Q(s', a')]$$

$$Q(s_t, a_t)$$

① learning

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha [R_t + \gamma \max_a Q(s_{t+1}, a)]$$

off-policy $Q(s_{t+1}, a)$

⇒ if we don't have probabilities we
have to do a large number of experiments
to find the ultimate
time and memory

Distribution model $\xrightarrow{\text{Experimenting}}$ Sample Model $\xrightarrow{\text{planning}}$ value

Distribution model $\xrightarrow{\text{planning}}$ values

Sample learning $\xrightarrow{\text{Value policy}}$

M	T	W	T	F	S
---	---	---	---	---	---

Date _____

Q-planning

loop forever:

select a state s

loop forever:

1. select s, a at random

2. send s, a to a sample model and get a sample next reward, R , and a sample next state, s'

3. Apply one-step tabular Q-learning to

$s, a \rightarrow R, s'$.

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

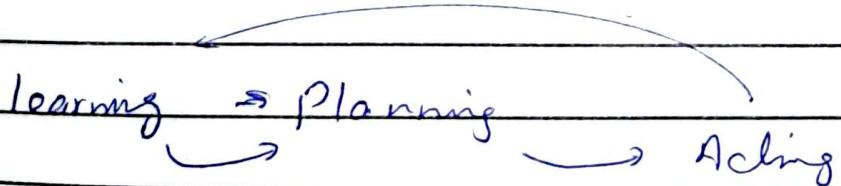
Sample Model

↓
Simulated behavior

Dyna Architecture

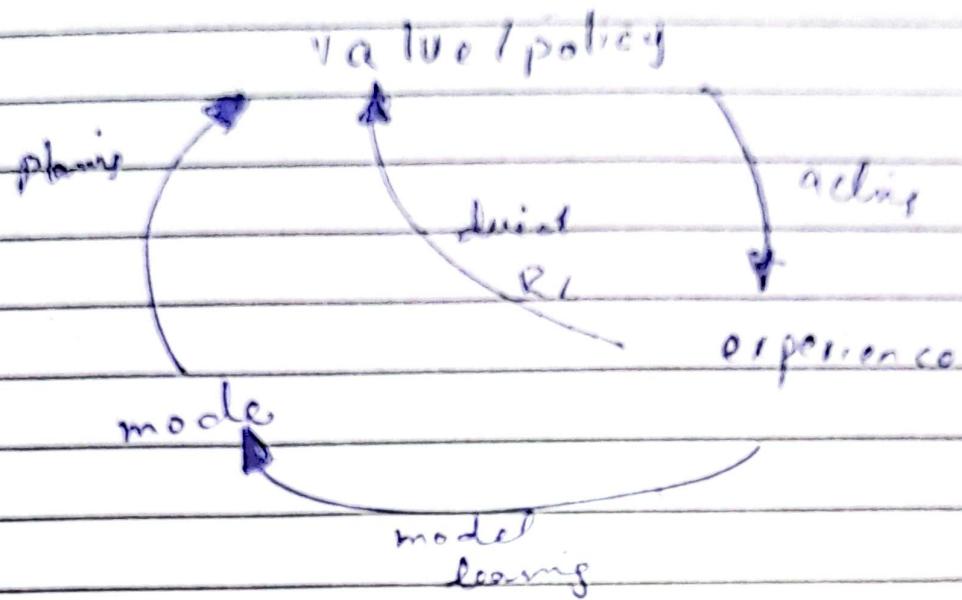
~~Dynamic~~
Dyna

→ learning → planning → Acting





MTWTFSS



Dyna Q \rightarrow for Planning and Learning using
Q-algorithm

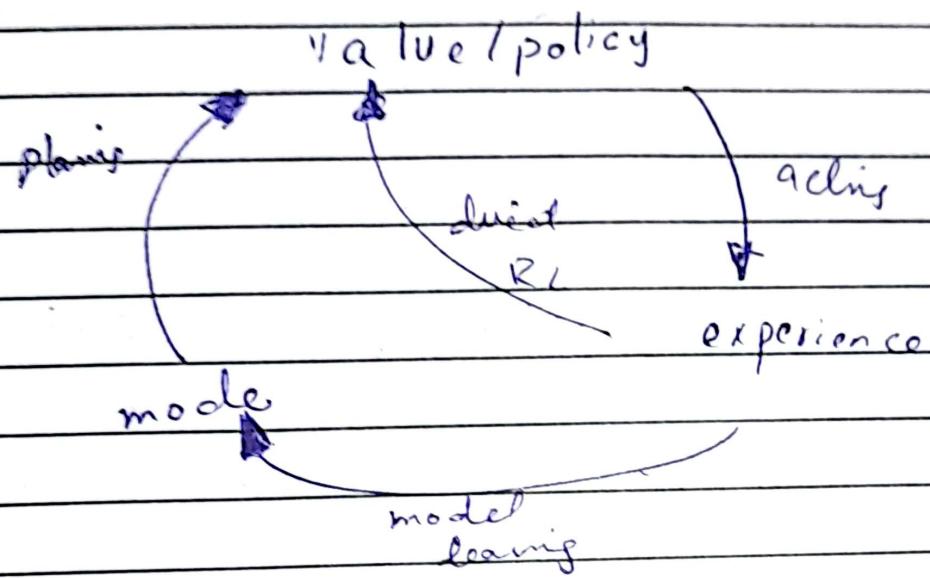
robot
standing on
+ cliff

Date _____

M T W T F S

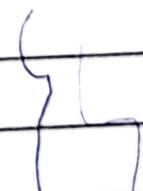


Castelli



Dyna Q \rightarrow for learning and planning using
Q algorithm

robot
standing on
a cliff



Date _____

MTWTF



Dyna Q

Dyna Q

action must be deterministic

Why the actions are deterministic in Dyna Q

If the action is ~~deterministic~~ stochastic
then the model will not be useful

For planning -

first-line action

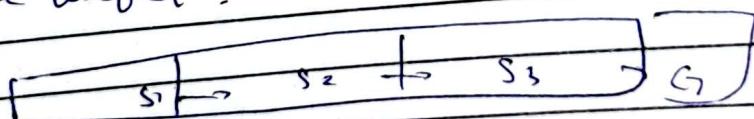
Model $(s, a,) \rightarrow s_2, o$

Stochastic

 $s_1, a_1 = s_2, 1$

Model Fail

not be useful.



S, A R, S'

S, A	R, S'
s_1, R_1	$0, s_2$
s_2, R_2	$0, s_3$
s_3, R_3	$1, G$

1st episode planning doesn't help you much

learning Q $(s_1, R_1) \leftarrow 0 + 1 \{ 0 + 0.9(0) \cdot 0 \}$

$$Q(s_1, R_1) \leftarrow 0$$

as

natively

all Q values
are 0Planning Repeat n times s_1, s_2

stop 0



M T W T F S

Date _____

In next step

$$Q(S_3, R) = 0 + 1 [1 + 0.9(0) \cdot 0] \\ = 1$$

$$G(S_3, R) = 1, G$$

Playing random choose S_2 S_2 update because of S_3

$$Q(S_2, R) = 0 + 1 [0 + 0.9(1) - 0] \\ = 0.9$$

Next time S_1

Problems

Model Incomplete

Model Inaccurate

→ If an action is inaccurate

so we update those actions reward larger so
that agent chooses them next time and
update their values.

We dynamically
update reward
based on how
long it is
since last request

Those state action who are updated after
long time and not change. we make their
reward larger so that ~~the~~ agent chooses
them

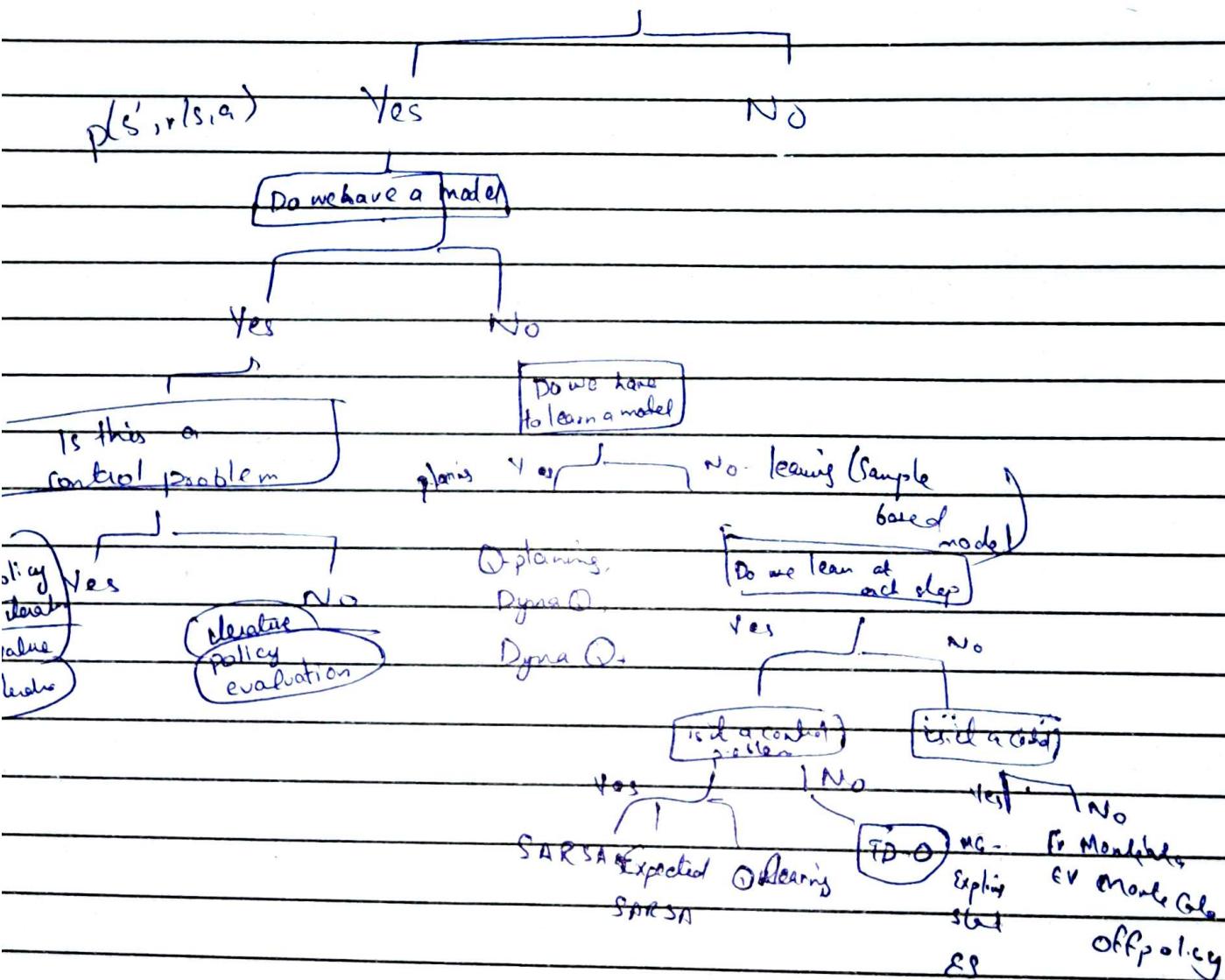
Dyna Q+

$$R(A) = R + \alpha \sqrt{\gamma} - \frac{\text{number of steps elapsed since the action } A \text{ has been taken}}{\text{hyperparameter}}$$

Date _____

MTWTF

Can we use tables to solve values



$$f(s) \rightarrow v$$

Date _____

MTWTF

Don't know the $v(s)$ True value $v_{\pi}(s) = y$ approximate value $\hat{v}(s, w) = y$

Need a function

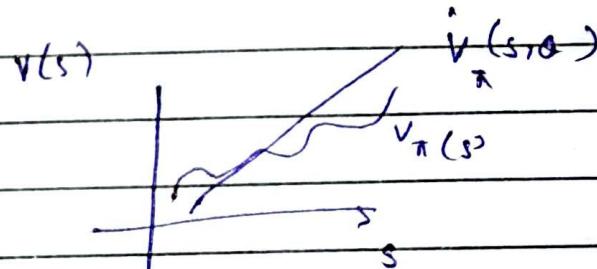
$$J(\theta) = \frac{1}{2} \sum_s (v_{\pi}(s) - \hat{v}(s, w))^2$$

In supervised

true values are given

but here we don't

know true values



if change
reduce one
error

when we change θ
error gets large for few steps
but gets small for few steps

We set θ to estimate
change accordingly

$$= \frac{1}{2} \sum_s u(s) (v_{\pi}(s) - \hat{v}(s, w))^2$$

Date _____

MTWTF

state which visit more ~~goes~~
 m. ↓
 more important
 Frequently visit states → ↓
 scale $u(s)$ value

$$\theta = \theta; - \frac{\partial J}{\partial \theta};$$

 ~~$u(s)$~~



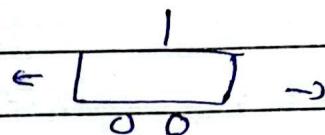
[M] [T] [W] [T] [F] [S]

Date _____

SARSA \rightarrow policy iteration

Q-learning \rightarrow value iteration

Cartpole



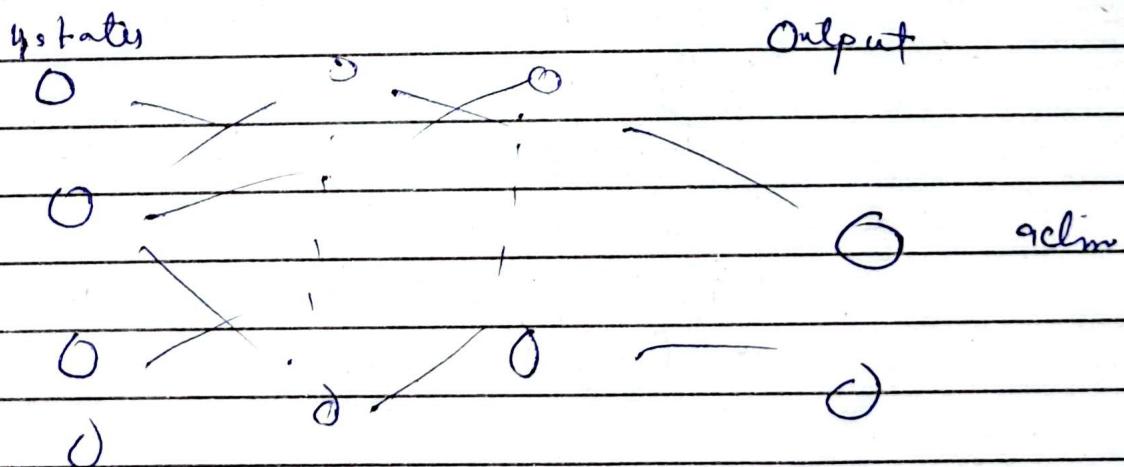
Actions 2 Left, right

State \rightarrow location, speed, angle, angular velocity

state space \rightarrow infinite

$q(s,a) \rightarrow$ can't store in a tabular data structure

Neural Network



forward pass

we have to change weights



Date _____

MTWTFSS

Loss function $J(\alpha) = \frac{1}{2} \sum (V(s) - V(s, a))^2$

reward Actual and Predicted value

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

predicted
values

\downarrow
Target values

$$\text{target} = R + \gamma \max_a Q(s', a')$$

From environment

From neural network

Algorithm:

1- Initialize:

- Neural Network (Q -network) to estimate $Q(s, a)$
- Replay buffer memory
- Hyperparameters like learning rate, epsilon, gamma

2- For each episode:

• Reset the environment to get an initial step.

• For each step,

1- Choose an action:

• Use ϵ -greedy strategy (explore or exploit)

2- Take the action

• Observe next state and reward

Date _____

(M) (T) (W) (T) (F) (S)



3 - Store the experience.

- Save (state, action, reward, next-state) to memory

4 Sample a batch from memory and:

- Compute Q-value estimates from the batch.

- Compute the Q-learning target:

$$\text{target} = r + \gamma \max_a Q(s, a)$$

- Compute loss and update network weights via backpropagation.

3 - Repeat until the agent learns to balance the pole.

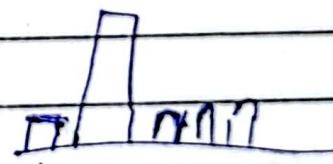
Policy Gradients

$$\pi(a|s)$$

$$Q(s, a)$$



ϵ -greedy



Choosing a_{ϵ}

Instead of action values

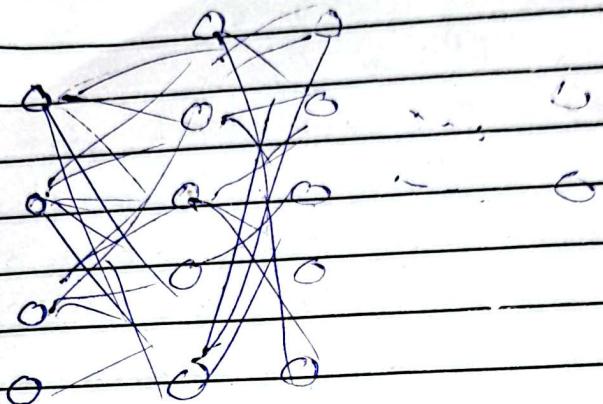
Get action probabilities
 a_1 one
 a_2 one
 a_3 one
 a_4 one

get action probabilities

We use softmax

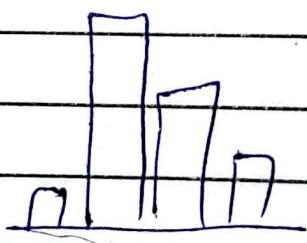
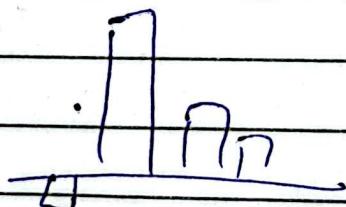
Date _____

(M T W T F S)



Softmax

Preference for actions

 $H(s, a)$ (relative)Adm value $\Delta(s, a)$ (actual values) $\Delta(s, a)$ $H(s, a)$ 

Preference

not heating
equally

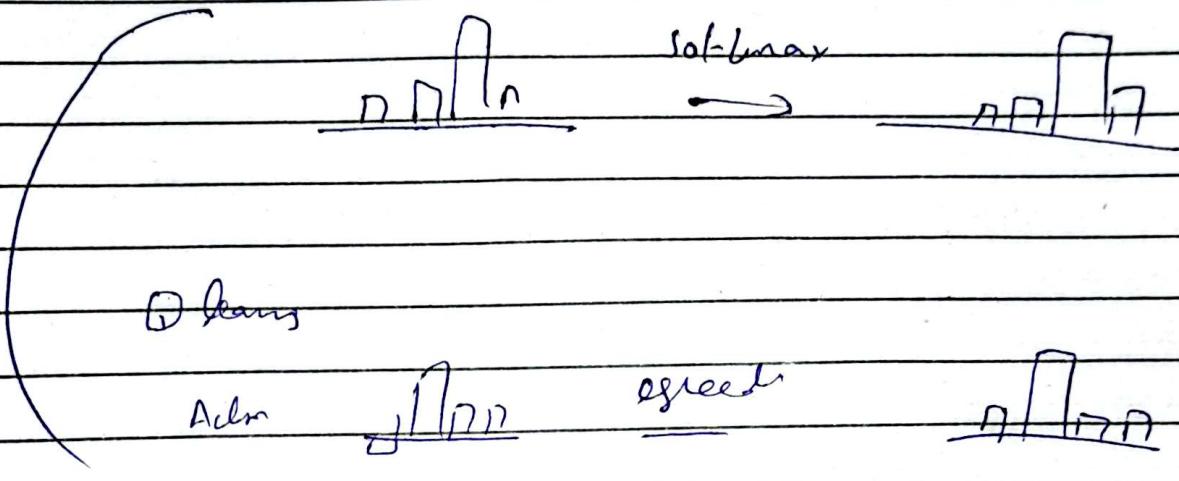
Date _____

MTWTFSS



In policy gradient we have
action preferences

$$\pi(a|s)$$



It consider the preferences.

Policy gradient use

Gradient Ascend

You maximize ~~step~~ the objective function

$$J(\theta) = r(\pi)$$

Average reward / action maximization

Average Return

$$\sum_a \pi(a|\theta) \sum_{s',r} p(s'|s,a) \cdot r = \sum_a \pi(a|\theta) \sum_s \sum_r \pi(a|s,\theta) \cdot r \cdot p(s',r|s,a) = \sum_s \sum_a \pi(a|s,\theta) \cdot r(s,a)$$

Date _____

MTWTFSS

$$\nabla \bar{r} = \sum_s u(s) \sum_a \pi(a|s, \theta) q_{\pi}(s, a)$$

↓
law

$$\theta = \theta + \nabla \bar{r}$$

$$\bar{r} = \sum \ln \pi(a|s, \theta) q_{\pi}(s, a)$$

↓

policy gradient using Monte Carlo

↓
called
PG IN FOLE

← TD

Actor
Critic