

Lecture 11 (Video 9)

M T W T F S



Date _____

Bellman's equation

state value function.

A function that tells us the value of a state
value of state tells us how good and bad an
state

value of state depends on the policy being followed

$$V_{\pi}(s) = \mathbb{E}(G_t | S_t = s)$$

Returns the discounted sum of reward.

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

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

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

$$\text{G}_{t+1} = G_{t+1, s} \mathbb{E}(G_{t+1} | S_{t+1} = s)$$

$$V_{\pi}(s) = \mathbb{E}(R_{t+1} + \gamma \mathbb{E}(G_{t+1} | S_{t+1} = s) | S_t = s)$$

$$V_{\pi}(s)$$

$$V_{\pi}(s) = \mathbb{E}(R_{t+1} + \gamma V_{\pi}(s') | S_t = s)$$

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

$$[\gamma + \gamma V_{\pi}(s')]$$

Date _____

MTWTF

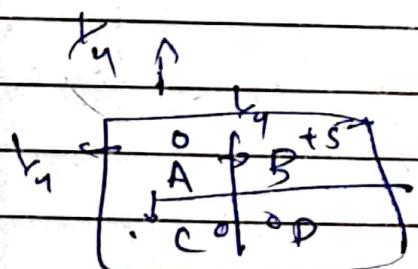
4 states

$$V_n(A) = 0.1 [0.1(5 + 0.9V_n(A)) +] \\ + 0.3 [0.2(3 + 0.9V_n(A)) + 0.4(2 + 0.9V_n(C)) + 0.3(1 + 0.9V_n(D))]$$

s	a	$\pi(a s)$	$+ 0.4 [$
A	left	$0.1 - \frac{A}{3}$	$0.1 [$
	right	$0.3 - \frac{C}{2}$	$0.2 [$
	bottom	$0.4 - \frac{B}{2}$	$0.4 [$
	up	$0.2 - \frac{D}{2}$	$0.3 [$

For deterministic

$$V_n(A) = \frac{1}{4} (0.9 V_n(C)) + \frac{1}{2} (0.9 V_n(A)) + \frac{1}{4} (5 + 0.9 V_n(B))$$



book example

Date _____

Action value

$$q_{\pi}(s, a) = \mathbb{E} (G_t | S_t = s, A_t = a)$$

$$= \mathbb{E} (R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a)$$

$$G_{t+1} = \mathbb{E} (G_{t+1} | S_{t+1} = s')$$

$$q_{\pi}(s, a) = \mathbb{E} (R_{t+1} + \gamma \mathbb{E} (G_{t+1} | S_{t+1} = s') | S_t = s, A_t = a)$$

$\overbrace{\quad\quad\quad\quad}$
 $V_{\pi}(s')$

$$q_{\pi}(s, a) = \mathbb{E} (R_{t+1} + \gamma V_{\pi}(s') | S_t = s, A_t = a)$$

~~$G(R_{t+1} + V_{\pi}(s'))$~~ = \mathbb{E}

$$V_{\pi}(s') = \mathbb{E} (G_{t+1} | S_{t+1} = s')$$

$$q_{\pi}(s, a) = \mathbb{E} (G_{t+1} | S_{t+1} = s', A_{t+1} = a)$$

$$V_{\pi}(s') = \sum_a \pi(a|s) \underbrace{\mathbb{E} (G_{t+1} | S_{t+1} = s', A_{t+1} = a)}_{\text{Eqn G. 9}}$$

$$q_{\pi}(s, a) = \mathbb{E} (R_{t+1} + \gamma \sum_a \pi(a|s) q_{\pi}(s', a))$$

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

Date _____

MTWTFSS

lecture 12

Video 10

Recap: Bellman Equations

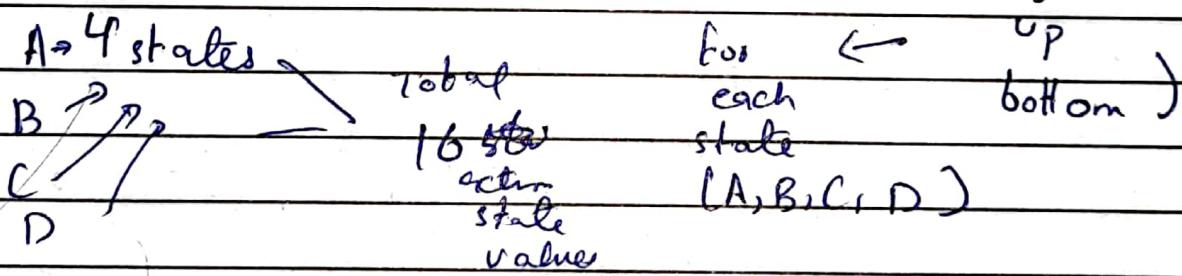
$$V_n(s) = \sum_a \pi(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma V_n(s')]$$

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

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

	\leftarrow	A	\rightarrow	B	\uparrow	S
C				D		

\Rightarrow Total 16 action equations
 4 actions
 left
 right



$$q_n(s', a')$$

Here $p(s', r|s, a)$ is ?

$$\begin{aligned}
 q(A, \text{left}) &= \frac{1}{4} [0 + 0.9 (\frac{1}{4} q(A, \text{left}) + \frac{1}{4} q(A, \text{right})) \\
 &\quad + \frac{1}{4} q(A, \text{up}) + \frac{1}{4} q(A, \text{down})]
 \end{aligned}$$

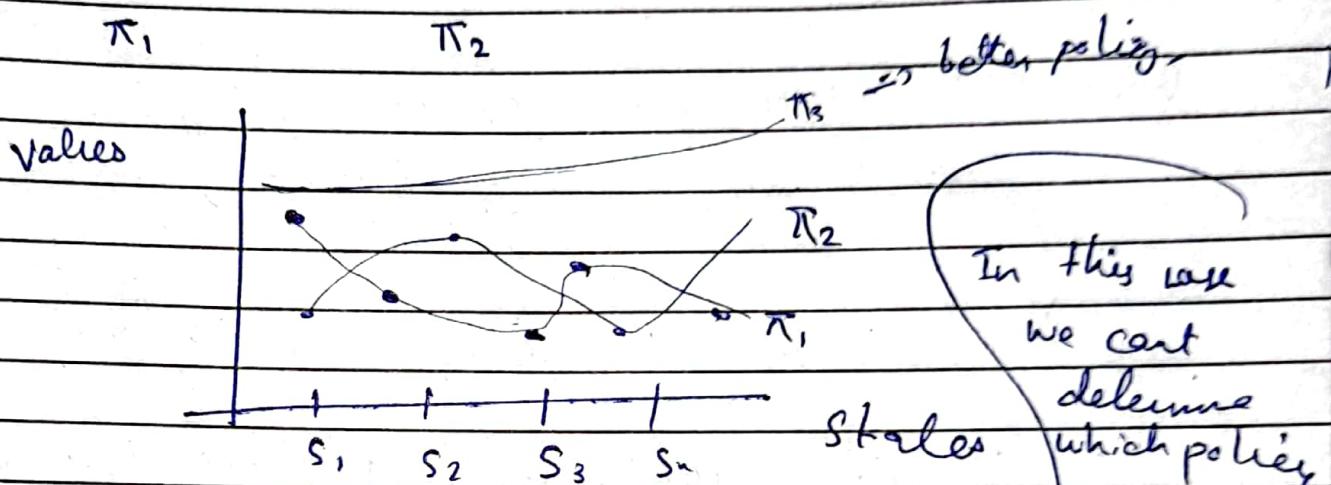
$$\begin{aligned}
 q(A, \text{up}) &= \frac{1}{4} \\
 q(A, \text{down}) &= \frac{1}{4}
 \end{aligned}$$





Now these action value comes under a policy which is random so these may be not be optimal values

When we say one policy is better than other value



$$\pi_1 \geq \pi_2$$

iff $V_{\pi_1}(s) > V_{\pi_2}(s)$ for each s belonging to S
 $\hookrightarrow S$

We can make a better policy by combining $\pi_1 \rightarrow \pi_2$

MOP

Control \rightarrow finding the best policy

Policy evaluation \rightarrow ~~finding~~ Comparing different policy for finding the best policy.

→ Can more than one optimal policy exist

Yes

$$V_\pi(s) = \max_{\pi} V_\pi(s)$$

Value of the state under optimal policy

$$q_\pi(s, a) = \max_{\pi} q_\pi(s, a)$$

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

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

What if two actions give you max value

↳ so we can choose from one of them
in this case even the deterministic policy become stochastic.

optimal policy can be stochastic?

Yes if there are more than one values/actions which give you max value.

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

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

Date

~~1, 2, 3, 4, 5, 6, 7, 8
9, 10, 11, 12, 13, 14, 15, 16~~


(M T W T F S)

$$V_n(s, a) \rightarrow q_n(s, a')$$

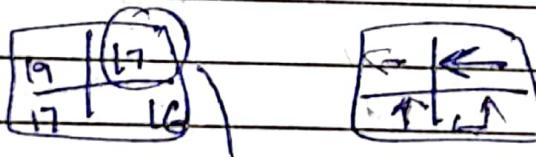
$$V_n(s') = \sum_a \pi(a|s') q_n(s', a')$$

$$V_{\pi}(s) = \sum_a \pi(a|s) q_{\pi}(s, a)$$

$$V_{\pi}(s) = \max_a q_{\pi}(s, a)$$

max function is not a linear function
 we can't solve it through system of linear equations

We solve it through DP.



V_n next action from the state

Optimal policy

$$\pi_n(s) = \arg \max_a q_n(s, a)$$

Date _____

MTWTF

Lecture 13

video 11

Revision

Optimal policy

$$V_\pi(s) = \max_\pi V_\pi(s) \quad s \in S$$

$$q_\pi(s, a) = \max_\pi q_\pi(s, a)$$

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

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

of each $V_\pi(s)$

20	18
13	15

$$V_\pi(0, 1)$$

$$q_\pi(s(0, 1), \text{up}) = 1(-1 + 0.9(18)) = A$$

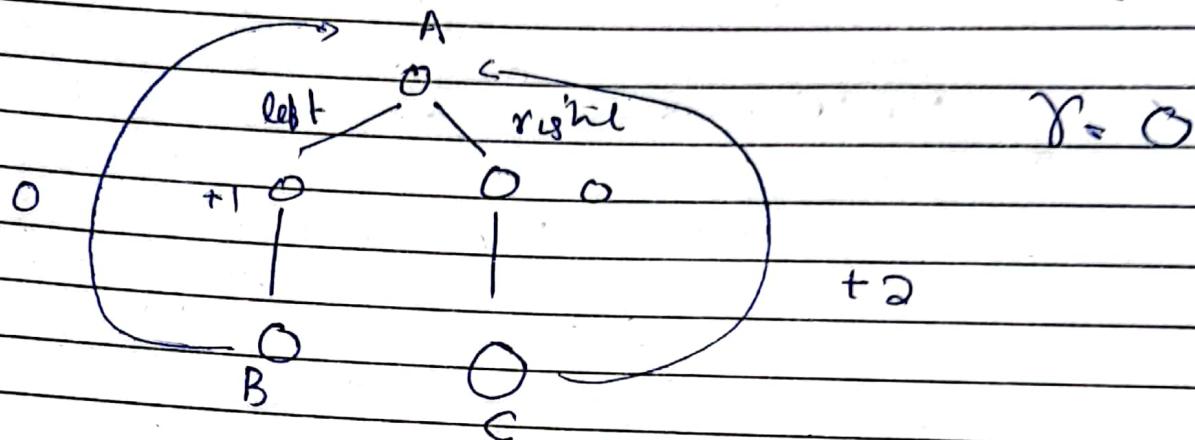
$$q_\pi(s(0, 1), \text{down})$$

$\frac{4}{11}$

which choose

which is max

$$V_{\pi_1}(s) \geq V_{\pi_2}(s) \quad \text{for every } s \in S$$



$$\mathbb{E}_{\pi}(s,a) \leq \mathbb{E}_{\pi'}(s,a) [r + \gamma \mathbb{P} V_{\pi'}(s')]$$

$$V(A) = 1[1 \cdot (1 + 0 \cdot V(B))] \\ V_{\text{left}}(A) = 1$$

which value is

$$V_{\text{right}}(A) = 0$$

greater $\frac{\text{left}}{\text{right}}$

As γ is 0 we are not giving any value to the future rewards that's why on $\gamma = 0$ we got $V(\text{left}(A))$ greater than $V(\text{right}(A))$

Now if we want to give value to future rewards

$$\gamma = 0.9$$

$$V_{\text{left}}(A) = 1[1 + 0.9(V(B))]$$

$$V_{\text{right}}(A) = 1[0 + 0.9(V(C))]$$

$$V(B) = 1[0 + 0.9(V(A))] \\ V(C) = 1[2 + 0.9(V(A))]$$

Date _____

(M) T W T F S

$$V_{\text{eff}}(n) = [1 + 0.9 \cdot (0.9^n V_A)]$$

$$V_L(n) = 1 + 0.81 (V_A)$$

$$\rightarrow V_L(n) \rightarrow$$

$$(1 - 0.81) V_A = 1$$

$$V_A = \frac{1}{1 - 0.81}$$

$$V(A) = \cancel{V_L} + \cancel{R_{k+1}}$$

$$V(A) = \sum_{n=0}^{\infty} \gamma^{k-t-1} R_{k+t+1}$$

$$V_{\text{eff}}(A) = 1 [1 + 0.9(0) + (0.9)^2(1) + (0.9)^3(0) + \dots]$$

 Geometric
Sequence

$$\sum_{n=0}^{\infty} r^n = \frac{1}{1-r}$$

$$= 1 \left[\sum_{n=0}^{\infty} (1)(0.9)^{2n} \right]$$

$$= \sum_{n=0}^{\infty} (1)(0.9^2)^n$$

$$V(A) = \frac{1}{1 - (0.9^2)}$$

$$V(A) = [0 + (0.9)(2) + (0.9)^2(0) + (0.9)^3(2) + \dots]$$

$$= \sum_{n=0}^{\infty} (1) 2 (0.9)^{2n+1}$$

$$= 2 \sum_{n=0}^{\infty} (0.9)(0.9)^{2n}$$

$$= 1.8 \sum_{n=0}^{\infty} (0.9^2)^n$$

$$= 1.8 \left(\frac{1}{1 - (0.9^2)} \right)$$

Date _____

MTWTFSS



Lecture # 14

Video 12

Dynamic Programming

MDP

$$\begin{aligned} V_\pi(s) \\ q_\pi(s, a) \end{aligned} \quad \left. \begin{array}{l} \text{evaluating} \\ \text{ } \end{array} \right\}$$

Goal \rightarrow find the optimal policy

Evaluation / prediction

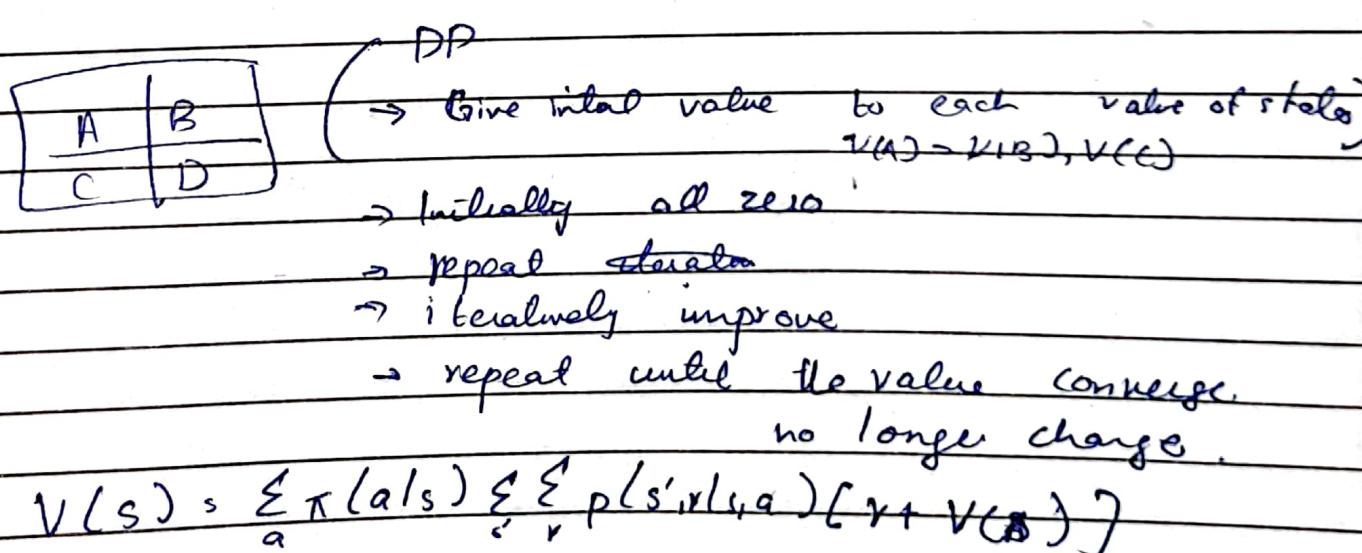
Control

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

Dynamic Programming

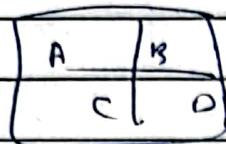
sub problem divide

store the solution and reuse it



Date _____

MTWTFSS



$$V(A) = \frac{1}{2} (0 + 0.9 V(A)^0) + \frac{1}{4} (0 + 0.9 (V(C)^0)) + \frac{1}{4} (5 + 0.9 (V(B)^0))$$

, solve it

$$V(A) = \frac{5}{4}$$

Find all

$$V(B), V(C), V(D)$$

repeat again putting these ^{new} value in equation

Book Example,

repe.

3	1	2	3
4	5	6	7
8	9	10	11
12	13	14	15

terminal state values
are 0

$$V(I) = \frac{1}{4} (-1 + 1(0)) + \frac{1}{4} (-1 + 0) + \dots$$

$$\Rightarrow -\frac{1}{4} - \frac{1}{4} + \frac{1}{4} = -\frac{1}{4}$$

$$V(I) = -1$$



Date _____

MTWTF

By

Lecture # 15

Video 13

 $k=0, \dots, N$

Dynamic Programming:

$$V_{\pi_k}(s) = \sum_a \pi_k(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma V_{k+1}(s')]$$

$$V_{\pi_k}(s) \rightarrow \sum_a \pi_k(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma V_{k+1}(s')]$$

$$\text{abs}(V_k(s) - V_{k+1}(s)) < \text{threshold}$$

Iterating policy evaluation \rightarrow Dynamic Programming
for finding $V(s)$

$$V_{\pi_1}(s) \geq V_{\pi_2}(s), \text{ for every } s \in S$$

$$\pi_1 \geq \pi_2$$

Initialize $V(s)$, for all $s \in S^+$

$$V(\text{terminal}) = 0 \quad \text{terminal state} = 0,$$

Loop:

$$\Delta \leftarrow 0$$

(Kitna difference arrha)

Loop for each $s \in S$:

$$v \leftarrow V(s)$$

$$V(s) \leftarrow \sum_a \pi_k(a|s) \sum_{s', r} p(s', r|s, a) [r + \gamma V(s')] \quad \text{Bellman's equation}$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|) \quad (25 \text{ states me se})$$

until $\Delta < 0$

8. ho ek ka previous iteration ki 25 states

se difference,

in 25 me se jo max value hogi wo Δ store hoga.)

We choose
max across all
the states

loop chalta rahega jis k time $\Delta < 0$ in hoga

Date _____

MTWTF

We use

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

if we use previous iteration $V(s)$ value to calculate new values in next iteration then we need 2 data structures one for storing old values and one for new,

$$V_{\text{old}}(s) \quad V_{\text{new}}(s)$$

but we can also use the values if we find them early during our iteration.
The only one data structure needed

Control:

How to find a better policy

Control → Policy iteration
Value iteration,

for deterministic policy

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

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

Are they same? No

$V(s)$ gives return when all actions are taken under the same policy

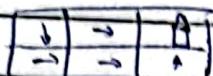
In $\pi(s,a)$ you first take action ~~and~~

~~the certain policy~~ ~~but then you~~

but

$q_{\pi}(s, a) \rightarrow$ first takes an action then follows a policy

In q we fix an action



↑ ↓

a is not according to this $\pi(s)$

$$q_{\pi}(s, \pi'(s)) \geq v_{\pi}(s)$$

~~if~~ \Rightarrow this action better

Policy Iteration

1. Initialization:

(random policy)

$v(s) \in \mathbb{R}$ and $\pi(s) \in A(s)$ arbitrarily for all $s \in S$

2. Policy Evaluation:

(value of stable funding)

loop:

$$\Delta \leftarrow 0$$

Loop for each $s \in S$:

$$v \leftarrow v(s)$$

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

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

until $\Delta < \theta$

Once value converges we improved

3. Policy Improvement:

policy-stable \leftarrow true

for each $s \in S$:

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

choose action which maximize the value

$$\pi(s) \leftarrow \arg \max \sum_{s', r, p} (s', r | s, a) [r + \gamma v(s')]$$

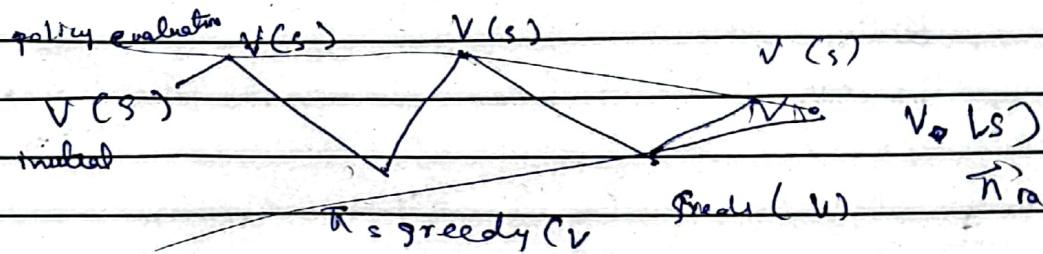
If old-action $\neq \pi(s)$, then policy-stable \leftarrow false

If policy-stable, then stop and return $V \approx v$ and $\pi \approx \pi_0$

else go to 2.

Date _____

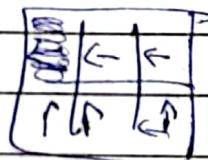
MTWTF



if you know reward and probability
 then your agent can
 find optimal pol. policy.

These are Model-based methods in which
 max action

0	14	-20
-18	-20	-20



lecture 16

video 14

Control - finding the optimal policy

Iterative Policy Evaluation

$\pi(s)$ → Prediction

Determining the value of state using DP

Policy Iteration

Iteratively finding the policy

Value Iteration

Policy iteration

Under a certain policy find

$V(s)$ on

the basis of old $\pi(s)$

using previous value of $V(s')$

threshold satisfy

Now

policy improvement

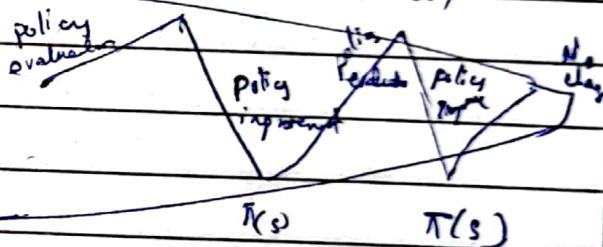
Now determine

(max action) to make $\pi(s)$ which

max $V(s)$ over policy

do until policy is stable

$V(s)$ $V(s')$



Bellman's equation

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

Bellman's optimality

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

* Policy converges faster than Value convergence

time wasted in policy iteration

so value iteration too is a solution.

We use Bellman's optimality equations in value iteration

Value iteration

Initialize $V(s)$ for all $s \in S$

Loop:

$\Delta \leftarrow 0$

loop for each $s \in S$:

$v \leftarrow V(s)$

$V(s) \leftarrow \max_a \sum_{s', r} p(s', r|s, a) [r + \gamma V(s')]$ (already optimal value we set after the loop finished)

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

until $\Delta < \theta$

Output a deterministic policy $\pi \approx \pi_{\pi^*}$

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



Date _____

MTWTFSS

In MDP \rightarrow we are doing planning (Information present)

Quick Recap:

MAB No time values
we estimate

MDP probabilities given
we find states

What if no Probabilities given

Sample-based Methods

We let the agent interact with environment

Monte Carlo Methods:

↳ randomness

is any experiment that has randomness and we tried to find some values by repeating the experiment a large number of times

Coin Example

A fair coin

$$P(\text{Head}) = 0.5 \quad P(\text{Tail}) = 0.5$$

If we don't know then we find ~~exp~~ probability on the basis of experimental

we find value of state on the basis of episodes run episodes multiple times

episode start & terminal state

$$V(s) = E(G_{t+1} | S_t = s)$$

↑ probability \times value

We call Average \Rightarrow Expected Value

Actual value $\xrightarrow{\text{experiment}}$

Average \approx Expected value (if ~~repeated~~ a large number of times)

Lecture 17 (Video 15')

Recap:

control - Policy Iteration

value iteration

Monte Carlo method

performing a large number of experiments for determining a value in a process that has randomness

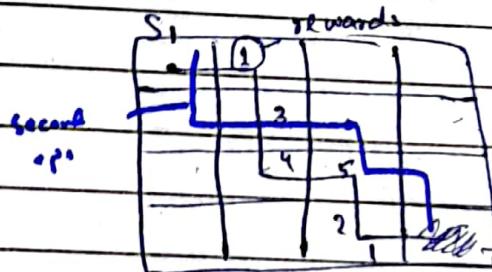
$$V(s) = \mathbb{E}(G_t | S_t = s)$$

if probabilities not given

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

we use average

$$G_t = R_{t+1} + rG_{t+1}$$

Experiment in RL: Episodes $t=0, \dots, t=N$, random policy

$E_1 \rightarrow S_0, A_0, R_1, S_1, A_1, R_2, S_2 \dots$

plus

columns of current state reward all rewards in future state till terminal

$$S_0 = 1 + 3 + 4 + 5 + 2 + 1 = 16$$

we start from terminal state to avoid repetition

Episode 1 $G(15) = 1$

$$G(14) = 2 + G(15) = 3$$

$$G(13) = 5 + G(14) = 8$$

$$G(12) = 4 + G(13) = 12$$

$$G(11) = 3 + G(12) = 15$$

$$G(10) = 1 + G(11) = 16$$

Episode 2

$$G(13) = 3$$

$$G(1) =$$

Episode 3

$$G(15) = 5$$

Two varieties of Monte Carlo

if within an episode a state is visited multiple times

$$G(15) = \frac{1+5}{2}$$

sum of state's value in all episodes

No. of episodes

- One approach to add the values (everytime it visited)
- one approach just use the first visited value

Date _____

MTWTF

First-visit MC

Only consider ^{the} state value $v(s)$
 when it is visited ~~but~~ in an episode.

Input a policy to be evaluated

Initialize:

$$V(s) \in \mathbb{R} \quad \text{for all } s \in S$$

$$\text{Returns}(s) \leftarrow \text{empty list} \quad \text{for all } s \in S$$

Loop forever (For each episode):

mean a
(large number
of times)

Generate an episode following $\pi: S_0, A_0, R_1, S_1, A_1, \dots, S_T, A_T, R_T$

$$G \rightarrow 0$$

step of

loop for each episode $t = T-1, T-2, \dots, 0$,

$$G \leftarrow rG + R_{t+1}$$

Unless S_t appears in S_0, S_1, \dots, S_{t-1} :Append G to Returns(S_t)

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

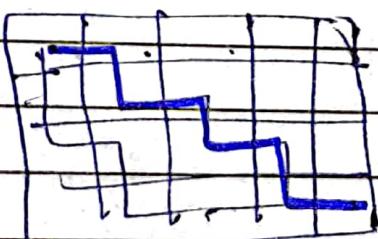
(Running average)

if probability is not given we

find values of state action values \rightarrow then find $\pi(s,a)$

$$V_\pi(s) = \mathbb{E}(G_t | S_{t+s})$$

$$q_\pi(s,a) = \mathbb{E}(G_t | S_{t+s}, A_{t+s} = a)$$

different episodes have k reasons \rightarrow policy stochastic \rightarrow policy deterministic $p(s',r)$ different, or reward different

$$\begin{aligned}
 \text{Episode 1: } Q(1, \text{down}) &= \frac{1}{25} \left[\text{On average take column 1} \right] \\
 Q(1, \text{right}) &= q_{\pi}(1, \text{right}) = 4
 \end{aligned}$$

Episode 2

if you did not take an action X in any episode
 but what if that was a good action,

Monte Carlo GS (Exploring Start,)

Initialize:

$$\pi(s) \in A(s) \quad \text{for all } s \in S$$

$$Q(s, a) \in R \quad \text{for all } s \in S, a \in A(s)$$

Returns $(s, a) \leftarrow$ empty list, for all $s \in S, a \in A(s)$

Loop forever (For each episode),

Choose $S_0 \in S$ & $A_0 \in A(S_0)$ randomly such that all pairs have probability 1/|A|

Generate an episode from S_0, A_0 following π

$$G \leftarrow \emptyset$$

loop for each step of episode, $t = T-1, T-2, \dots, 0$:

$$G_t \leftarrow r G + R_{t+1}$$

Unless the pair S_t, A_t appear in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$
 Append G_t to Returns (S_t, A_t)

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

$$\pi(S_t) \leftarrow \text{argmax}_a Q(S_t, a)$$

Exploring Start)

has episode to explore new state or get in
 we start knew

has state or hr action to get return
 value average)



Date _____

MTWTF



Lecture 18 (video (G))

Monte Carlo Methods (Sample-based Method)

- ↳ policy evaluation
- ↳ control

Policy Evaluation $\rightarrow \pi$ (fix a policy)

Finite no. steps

$$V_{\pi}(s)$$

Episodes: $S_0, A_0, R_1, S_1, A_1, \dots$ Episodes $\rightarrow (1) 2 3 4 \dots N$

S_i, A_i, \dots States

$$S_1 \quad \text{Return}(S_1) \quad \text{Return}(S_1)$$

$$S_2 \quad \text{Return}(S_2)$$

$$S_3 \quad \text{Return}(S_3) \quad \text{Return}(S_3)$$

$$S_4 \quad \text{Return}(S_4)$$

$$\text{Avg(Return)} = V_{\pi}(s)$$

Two Types of Monte Carlo \rightarrow First-visit MC

\rightarrow Every-visit MC

$$\text{Running average} \rightarrow \text{sample avg} \quad Q(a) = \frac{1}{n} \sum R_i$$

Incremental average

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

In previous module
MDP

We can find optimal value

$$\pi^* = \arg \max \left(\mathbb{E}_{\pi^*} [\text{Return}(s, a)] \right)$$

We can find control there
because $p(s_{t+1}, a)$ but

in Monte Carlo we can't
as there is no $p(s_{t+1}, a)$ so

for finding optimal value we
take $Q(s, a)$

$$V_{\pi}(s) = V(s) + \frac{1}{n} (\text{Return} - V(s))$$

$$V(s) = V(s) + \frac{1}{n} (G_t - V(s))$$

Date _____

MTWTFSS



Episodes 1, 2 3 4

s_0, a_0 Return()

s_0, a_1

s_0, a_2

s_1, a_3 Return()

s_1, a_1

s_2, a_2 Return()

So there was the problem that some state, action pair will not be ~~visited~~ visited

Solution:

In episode start we randomly choose ~~the~~ pair whose s, a value probability is positive at least choose every (s, a) value pair one time but it is not enough

Monte Carlo (ES)

so we want to do exploration here

greedy approach

→ here it



In MDP greedy actn is not good because there we have information about environment, probabilities are given

is not good

good

Epsilon - greedy

ϵ -greedy

With ϵ , small value $> 0 \rightarrow$ choose random actn with $1 - \epsilon$ choose greedy actn

Date _____

MTWTF

 A_1, A_2, A_3, A_4 $P(A_3) = \text{max}$

best

 $\varepsilon \rightarrow \text{random action}$

$$\varepsilon = P(A_1) \text{ or } P(A_2) \text{ or } P(A_3) \text{ or } P(A_4)$$

$$\varepsilon = P(A_1) + P(A_2) + P(A_3) + P(A_4)$$

$$\varepsilon = 4P(A_i)$$

$$\begin{aligned} P(A_1) &= P(A_2) \\ &= P(A_3) \\ &= P(A_4) \end{aligned}$$

$$P(A_i) = \frac{\varepsilon}{4} \quad (\text{choosing a random action probabilities})$$

 $|A|$ no of actions

size of action space

$$P(A) = \frac{\varepsilon}{|A|}, \quad \frac{\varepsilon}{4} \text{ is 1 action}$$

$$P(A_1) = \frac{\varepsilon}{|A|} = P(A_2) \quad \text{probability of choosing a random action}$$

Probability of choosing a greedy action

(As greedy action can be chosen)

$$P(A_3) = 1 - \varepsilon + \frac{\varepsilon}{|A|}$$

while exploration

$$P(A_1), P(A_2), P(A_4) \leq \frac{\varepsilon}{|A|}$$

In exploring start with of greedy we ~~choose~~ use Epsilon greedy



On-policy first MC control (for ϵ -soft policies)

small $\epsilon > 0$

Initialize:

$\pi \leftarrow$ an ϵ -soft policy

$Q(s, a) \in \mathbb{R}$

Returns (s, a)

In start no policy
an ϵ -soft policy

In the end it
becomes epsilon-greedy

Repeat forever (for each episode):

Generate an episode following π : $S_0, A_0, P, \dots, S_{T-1}, A_T, R_T$

$G \leftarrow 0$

Loop for each step of episode $t = T-1, T-2, \dots, 0$:

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

Unless the pair (S_t, A_t) appears in $S_0, A_0, S_1, A_1, \dots, S_{t-1}, A_{t-1}$

Append G to Returns(S_t, A_t)

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

$A^* \leftarrow \arg\max_a Q(S_t, a)$

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

$\pi(a|S_t) \leftarrow \begin{cases} \text{Greedy} & \text{if } a = A^* \\ 1 - \epsilon + \epsilon / |A(S_t)| & \text{if } a \neq A^* \\ \epsilon / |A(S_t)| & \text{Non greedy if } a \neq A^* \end{cases}$

Epsilon-soft.

superset of epsilon-greedy

Every epsilon-greedy is epsilon-soft but

not every epsilon-soft is epsilon-greedy.

Lecture 19 (Video 17)

Monte Carlo - Sample based
 $Q(s, a)$

$$\pi = \text{argmax}_a Q(s, a)$$

- ① Exploring starts - episode starts in
 • State, action
 • Exploit explore
 how to

② Epsilon-greedy

$$\begin{cases} \epsilon / |A(s_t)| & \epsilon \rightarrow \text{non-greedy actn} \\ 1 - \epsilon / |A(s_t)| & 1 - \epsilon \rightarrow \text{greedy action} \end{cases}$$

$$Q_{t+1}(s, a) = Q_{t+1}(s, a) + \text{stepsize} (G_{t+1} - Q(s, a))$$

$$= Q_{t+1}(s, a) + \text{stepsize} ((P_{t+1} G_{t+1}) - Q(s, a))$$

Epsilon-soft

$$P(A) \geq \frac{\epsilon}{|A|}$$

Has action a_i choose a_i with probability π
 (at least $\frac{\epsilon}{|A|}$) no $\epsilon = 0.1, |A| = 4$

$$P(A_1) \geq 0.1$$

$$P(A_2) \geq 0.1$$

$$P(A_3) \geq 0.1$$

$$P(A_4) \geq 0.1$$

Epsilon-greedy

No-greedy

$$P(A) = \frac{\epsilon}{|A|}$$

Greedy

$$P(A) = 1 - \epsilon + \frac{\epsilon}{|A|}$$

Equi-probable random policy is also
epsilon soft

has action k
chosen here
k equal
chances here

$$P(A) = \frac{1}{|A|} \rightarrow \epsilon \text{ s.t.}$$

Epsilon - greedy probability

In previous algo

$$\text{greedy} \rightarrow \pi = \arg\max_{\pi} Q(s, a)$$

Being greedy is being optimal as
it may be wrong so for exploration we
use epsilon-greedy
but greedy action maybe sometime a good
action

But epsilon soft is stopping you from choosing
greedy action even in the case when greedy
action is a correct
so epsilon-soft is also not optimal

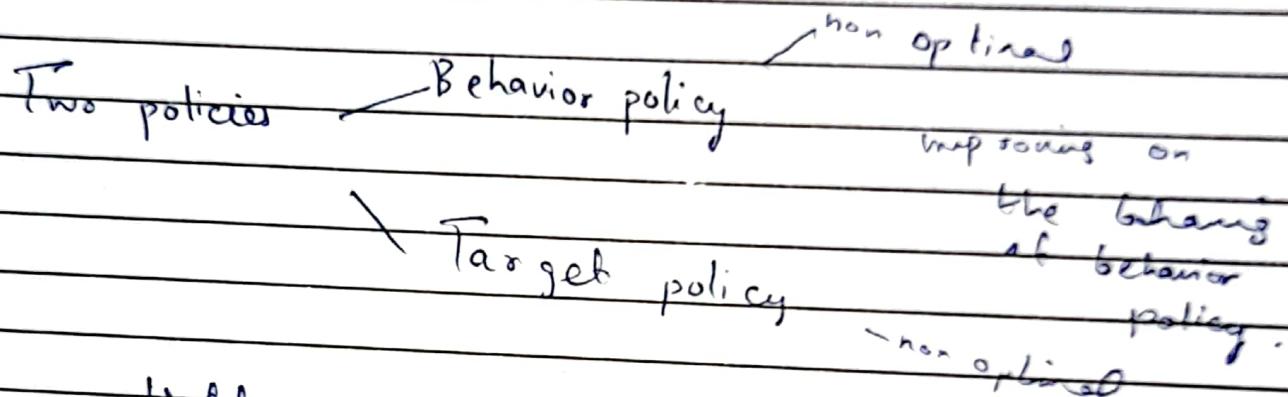
both greedy and epsilon-soft are
not optimal but if we compare which is
better epsilon soft \rightarrow greedy
So we want to improve
epsilon soft



Date _____

MTWTFSS

if an agent is learning from its behavior then it may not be a good thing if the behavior is not optimal



Until now the algorithms we see they both are equal

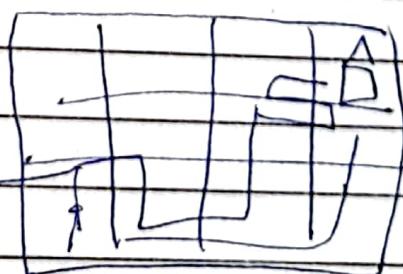
On policy \rightarrow Behavior policy = Target policy

if behavior policy is not optimal then target is not optimal

How to separate both these

off-policy

behavioral policy



k through k one

or learn during