

## Monte Carlo R-2.

Planning by D.P. can be done when MDP is known.

It requires knowledge of transition probabilities  $P(s, a, s')$

But when MDP is unknown, model free methods need to be used

Monte carlo is a model free method that learns directly from episodes.

MC learns from full episodes. It does not use bootstrapping.

All episodes must terminate

MC use empirical mean return instead of expected return.

Goal: Learn  $V_s$  from episodes of experience under policy  $\pi$

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

But in monte carlo, transition probability is not available.

Hence there is no weights of probability to multiply the rewards by.

Hence we cannot use  $E_{\pi}$  (expectation of reward of policy)

In monte carlo, empirical transition  $E$  is used

Value function for a policy is calculated by running the episodes and calculating reward for every state

$$\begin{array}{c} V \\ \nearrow \quad \searrow \\ \pi \end{array} \rightarrow V_{\pi}^*$$

Iteratively the policy is also improved.

# First Visit Monte Carlo Policy Evaluation

To evaluate  $V$  of a state  $s$

The first time that state  $s$  is visited in an episode, (time  $t_i$ )

Increment counter  $N(s) \leftarrow N(s) + 1$   
Increment total return  $S(s) \leftarrow S(s) + G_{t_i}$

At end of all episodes

$$V(s) = \frac{S(s)}{N(s)}$$

Here  $G_{t_i}$  is total discounted reward obtained from  $t_i$  time  $t_i$  till end of episode

$$G_{t_i} = \sum_{n=i}^{\text{end}} \gamma^{n-i} R_n = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots R_T$$

Ignore subsequent visits to the same state

The visits are done w.r.t. policy

---

## Every Visit MC

Do updation every time state is visited

First Visit MC	Every Visit MC
Unbiased estimate of value function	Biased because the returns for subsequent visits are not independent.
Higher Variance	Lower variance since it utilizes all the data

# Variance in Monte Carlo

In Monte Carlo, the calculations for all the states is not done unlike D.P.

Each episode is influenced by random events, stochastic environments or different consequences of actions can be random.

The total return  $G$  can vary significantly from one episode to next, even if the agent visits the same state multiple times.

In such cases, this causes a lot of variance in the model.

Another source of variance is that MC waits until the episode is over to update values, so later rewards that are not directly connected to the current state are also considered.

Unusual rewards at any point in the episode can impact the value estimate for a state, resulting in higher variance.

For example, episode ends with a rare large reward, it can lead to a significant jump in the value update causing large swings in learning.

$V_\pi$  in Monte Carlo

$V$  in M.C. is always with respect to  $\pi$

This is because reward of all future states is considered in  $G_t$

$V$  is the goodness of states when agent follows  $\pi$

the future states visited depends on  $\pi$

$$V(s) \leftarrow V(s) + \frac{1}{N} (G_t - V(s))$$



$V$  in MC depends on the policy taken  $\pi$ .

$$r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^n r_{t+n}$$



changing  $\pi$  will change  $V$ , since  $\pi$  determines the actions that the agent will take from state  $s$ .

depends on policy  $\pi$   
 $\therefore V(s)$  also depends on  $\pi$

eg agent takes first move in chess as  $E_1$  as per its policy  $\pi$

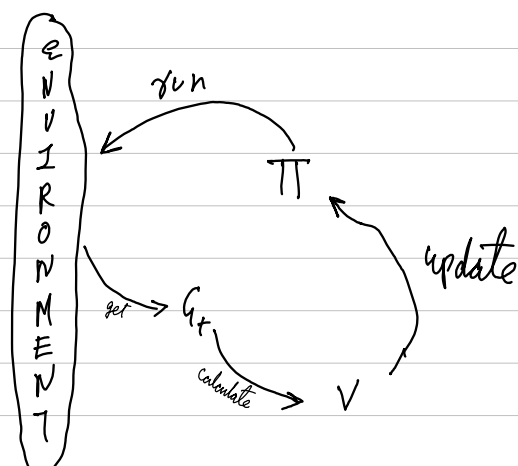
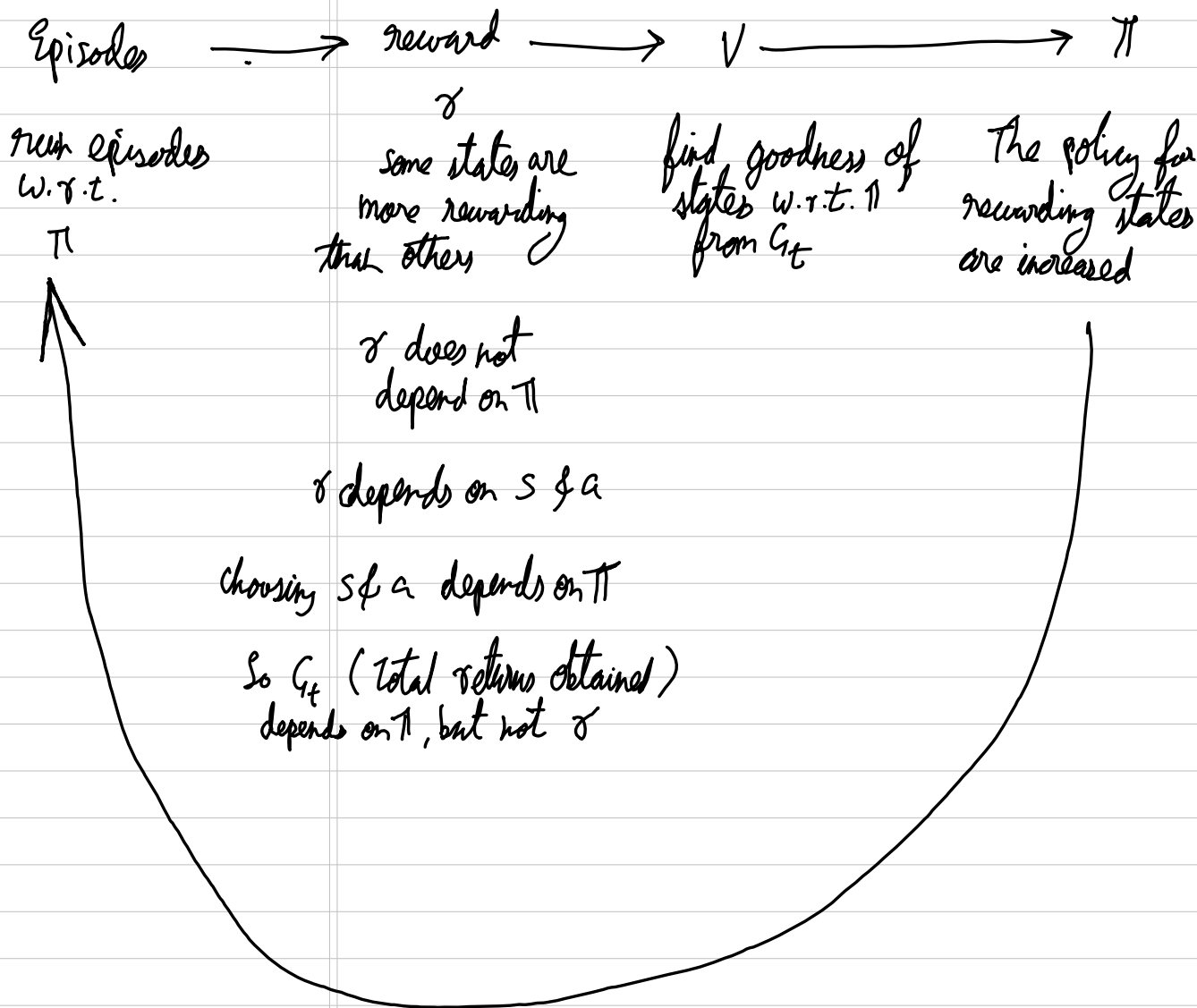
Case 1: agent is not trained properly and converges to non optimal  $\pi$

$V$  of  $E_1$  is less

Case 2: agent is trained properly and can win nearly every game

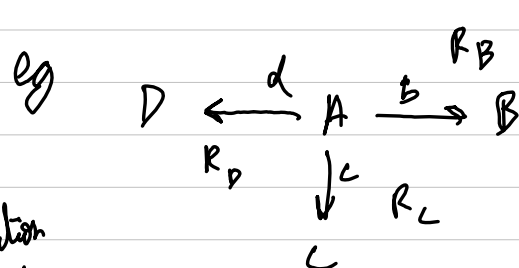
$V$  of  $E_1$  is high

The actual value of  $V(E_1)$  is not determined but the value functions determined are under policy  $\pi$ .



theoretically,  $G_t$  will be

$G_t$  run episodes  $\rightarrow$  Probability of action w.r.t.  $\pi$   $\times$  reward obtained



in iteration over many episodes

$$\therefore G_{t+A} \approx b \times R_B + c \times R_C + d \times R_D$$

$G_t$  is a measure how much rewards will be obtained on following policy  $\pi$

$$\pi = \begin{bmatrix} b \\ c \\ d \end{bmatrix}$$

$$V_A \leftarrow V_A + \alpha (G_{t+A} - V_A)$$

$\therefore V_A$  depends on  $[b, c, d]$  i.e.  $\pi$

Then improve  $\pi$  greedily wrt  $V$

Note  $\rightarrow$  states explored BCD also depends upon  $\pi$

Actions explored bcd also depends upon  $\pi$

rewards  $\gamma$  don't depend on anything

## Incremental MC Updates (Prediction)

Update  $V(s)$  incrementally after each episode instead of doing at the end

$$N(s_t) \leftarrow N(s_t) + 1$$

$$V(s_t) \leftarrow V(s_t) + \frac{1}{N(s_t)} (G_t - V(s_t))$$

This is exact same as the MC updates

However, we can now modify this equation to give importance to new episodes.

$$V(s_t) \leftarrow V(s_t) + \alpha (G_t - V(s_t))$$

This is done in order to forget old episodes.

Useful for modelling environments that are non stationary

## $\epsilon$ - greedy Policy

In  $\epsilon$ -greedy policy, we make choice between exploration & exploitation

While running the episode

Probability of  $\epsilon \rightarrow$  chose random action  
(explore)

$1-\epsilon \rightarrow$  chose highest estimated  
value (exploit)

So agent will exploit with prob  $(1-\epsilon)$  &  
explore with probability  $\epsilon$

# MC for deterministic vs Non deterministic

## Deterministic

$$\pi(s) = a$$

That means episodes will choose only the best so far case while exploring

$$\pi(s) \leftarrow \underset{a}{\operatorname{Argmax}} Q(s, a)$$

$$Q(s, a) \leftarrow Q(s, a) + \alpha (r_t - Q(s, a))$$

## Stochastic policy

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

That means the agent is allowed to explore even in the exploitation phase

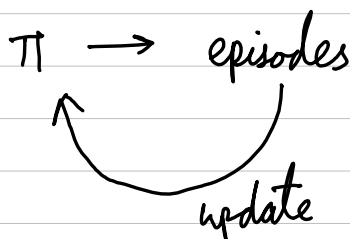
$$\pi(a|s) = \frac{e^{\frac{Q(s, a)}{\tau}}}{\sum_b e^{\frac{Q(s, b)}{\tau}}}$$

where  $\tau$  is temperature parameter that controls the level of exploration



## On Policy Monte Carlo

The methods till now were on policy, that is we used  $\pi$  to generate episodes



Policy to be optimized  $\pi$ ,  
same policy is used for generating episodes

But this is a compromise

It aims to learn action values  $V$  or  $Q$  that are dependant on subsequent optimal behaviour  $\pi^*$  but have to behave suboptimally to explore all actions and find optimal action

Learns  $V, Q$  not for  $\pi^*$ , but for a

"near optimal" policy that still explores  $\pi^\epsilon$   
eg epsilon greedy

Remember  $V \& Q$  are always w.r.t.  $\pi$

so we are actually finding  $V^{\pi^\epsilon}$  & not  $V^{\pi^*}$

Off policy MC is used to solve this

# off Policy Monte Carlo

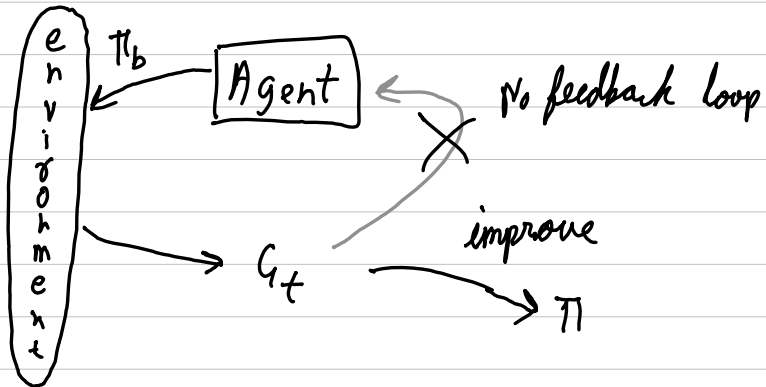
We use 2 separate policies  $\pi$  &  $\pi_b$

$\pi \rightarrow$  target policy

$\pi_b \rightarrow$  behavioural policy

Optimize  $\pi$  while generating behaviours from  $\pi_b$

$\pi_b$  can explore  
 $\pi_b$  can even follow some heuristic



But now how can we find  $V^\pi, Q^\pi$  & not  $V^{\pi_b}, Q^{\pi_b}$ ?

for this, we need importance sampling

# Importance Sampling

Importance sampling is used to evaluate properties of  $\pi$  when samples are generated by different distribution  $\pi_b$

When we generate episodes from  $\pi_b$  & get return  $G_t$ , we need to normalize the return.

$G_t$  : How much rewards are obtained when model follows  $\pi_b$

$G'_t$  : How much rewards would have had been obtained if model had followed  $\pi$

We need  $G'_t$  & not  $G_t$

$$G'_t = \frac{\text{likelihood of } \pi \text{ taking action}}{\text{likelihood of } \pi_b \text{ taking action}} \times G_t$$

$$= p_{0:T} \times G_t$$

where  $p_{0:T} = \prod_{t=0}^{T-1} \frac{\pi(a_t | s_t)}{\pi_b(a_t | s_t)}$  (Importance sampling ratio)

Annotations:   
 - "Action that caused the reward" points to  $a_t$  in the numerator.   
 - "action taken by  $\pi$  at time  $t$ " points to  $\pi(a_t | s_t)$ .   
 - "state at  $t$ " points to  $s_t$  in the denominator.   
 - "action taken by  $\pi_b$  at time  $t$ " points to  $\pi_b(a_t | s_t)$ .

$$\therefore V(s) = \frac{\sum_{t=0}^{T-1} p_{t:T-1} \times G_t}{n}$$

Annotations:   
 - "t is states through time" points to  $t$  in the sum.   
 - "no of times  $s$  is visited in all eps." points to  $n$ .

This simply is going through each episode where  $s$  is visited, then summing returns recieved after  $s$  was visited but weighting them by importance sampling ratio. finally taking the average } ordinary importance sampling

Here there is a variation (weighted importance sampling) that is preferred

$$V(s) = \frac{\sum_{t=0}^{T-1} p_{t:T-1} \times G_t}{\sum_{t=0}^{T-1} p_{t:T-1}}$$

The advantage is that we can learn from any historic data, given data  $G_t$  &  $\pi_b$

How much rewards are obtained on following  $\pi_b$



How much rewards would have had been obtained on following  $\pi$



Update  $V$



Update  $\pi$

Importance sampling is a variance reduction technique

# Monte Carlo Control

In MC control, the goal is not only to evaluate a fixed policy but also to improve it wrt actions

Policy improvement is done greedily

$$\pi'(s) = \arg \max_a q^{\pi}(s, a)$$

In cases where  $P(s'|s, a)$  is known

$$\pi'(s) = \arg \max_a \sum_{s'} P(s'|s, a) V^{\pi}(s')$$

In MC prediction, goodness of states is taken (state, value)

In MC control, goodness of state & actions are considered (state, action)

Prediction	Control
Policy is supplied and goal is to check how well it performs	Control is where the policy is not fixed and goal is to find the optimal policy

Prediction problems can be used for control when transition probabilities  $P(s, a, s')$  are known

Useful in training Games, navigation

## Estimating Action value function

Similarly to estimate action value function  $Q^\pi(s, a)$  let  $D(s, a)$  be set of all time steps  $t$  at which the state action pair  $(s, a)$  is visited.

Then the monte carlo estimate of  $Q^\pi(s, a)$

$$Q^\pi(s, a) = \frac{1}{|D(s, a)|} \sum_{t \in D(s, a)} G_t$$

for Incremental case

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha (G_t - Q(s_t, a_t))$$

$$\pi(s) = \underset{a}{\operatorname{argmax}} Q(s, a)$$