

Info state | markov state

$$P(s_{t+1} | s_t) = P(s_{t+1} | s_1, s_2, \dots, s_t)$$

future is independent of the past given the present

State representation defines what happens next

Fully observable environments

agent directly observes env state

$$o_t = s_t^a = s_t^e$$

Agent state = env state = info state

Markov  
Decision  
process  
(MDP)

Partially observable env

agent indirectly observes env state

e.g.: robot with a camera w/o GPS

- : trading agent with only current prices
- : poker agent

agent state  $\neq$  env state

agent must make their own ref.  $s_t^a$

Partially  
Observable  
MDP  
(POMDP)

different ways:

- complete history

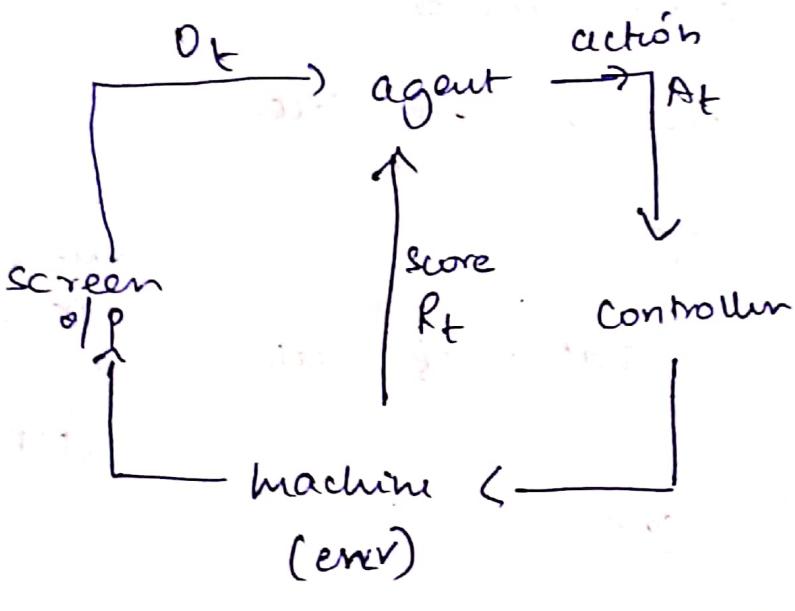
$$s_t^q \in H_t$$

- Beliefs of env state

$$s_t^a = (P[s_t^e = s^1], \dots, P[s_t^e = s^n])$$

- RNNs

Atari eg:



(1024 bit env state)

## Inside an RL agent

### main concepts

- Policy : how it picks action
- value : how good is it in a particular state
- Model : agent rep of the env

### Policy

A map from state to action

- Deterministic policy

$$a = \pi(s)$$

- Stochastic policy

$$\pi(a|s) = P(A|S=s)$$

### Value

- Prediction of expected future reward

$$V_{\pi}(s) = E_{\pi}[R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots | S_t = s]$$

↳ future rewards

↷ → discounting

## Modul

- predicts what env will do

- transition model

- predicting next state

- reward model

- predicting next reward

$$P_{ss'}^a \geq P[S' = s' | S = s, A = a]$$

$$R_s^a = \mathbb{E}[R | S = s, A = a]$$

Characterizing RL agent

(1)

- Value based

• No policy (implicit)

• Value function

- Policy based

• Policy

• No value fu

- Actor critic

- Policy

- Value function

## (2) - Model free

- Policy and/or Value fn
- No Model

## - Model based

- Policy and/or Value fn
- ~~No Model~~

## Problems in RL

---

### - RL

- env is unknown (robots)
- env interacts with env
- env improves its policy

### - Planning

- Rules known
- env model known
- agent performs computation
- improves policy
- could plan ahead using methods like tree search.

## Exploration & exploitation

- trial & error
- goal is to discover good policy.
- exploration - more info on env  
exploitation - exploit available info to get more reward

## Markov Decision Processes

- formally describes an environment for RL
- Almost all RL problems can be formalized as MDP
- Any partially observable problem can be converted into a MDP
- Bandits are MDPs with one state

## State transition matrix

$$\text{transition probability} = P[S_{t+1} = s' | S_t = s]$$

## Markov process

sequence of random states that have Markov property.

A Markov Process is a tuple  $(S, P)$

$S \rightarrow$  set of states

$P \rightarrow$  transition probability matrix

## Markov Reward Process

$(S, P, R, \gamma)$

$R \rightarrow$  Reward fn

(immediate reward)  $R_s = E[R_{t+1} | S_t = s]$

Return  $G_t \rightarrow$  total discounted reward from time  $t$

$$G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots = \sum_0^{\infty} \gamma^k R_{t+k+1}$$

discount is the present value of future rewards

## Why discount?

- Because we don't have a perfect model. So there will be uncertainty in the future.
- Mathematically convenient
- Avoids infinite returns

- On financial setting money now might be "better" than delayed money.
- Animal / human behavior.

Value function

long-term value of state  $s$

$$v(s) = E[g_t | s_t = s]$$

Bellman Eq

$$v(s) = E[g_t | s_t = s]$$

$$= E[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | s_t = s]$$

$$= E[R_{t+1} + \gamma [R_{t+2} + \gamma R_{t+3} + \dots] | s_t = s]$$

$$= E[R_{t+1} + v(s+1) | s_t = s]$$

$$= E[R_{t+1} + v(s+1) | s_t = s]$$

$$V_2 = R + \gamma P_V$$

immediate reward

$$\begin{bmatrix} V_1 \\ V_2 \\ \vdots \\ V_n \end{bmatrix}$$

$$\begin{bmatrix} R_1 \\ R_2 \\ \vdots \\ R_n \end{bmatrix}$$

$$+ \gamma$$

Transition matrix

$$\begin{bmatrix} P_{11} & \dots & P_{1n} \\ \vdots & \ddots & \vdots \\ P_{n1} & \dots & P_{nn} \end{bmatrix}$$

value of end cap

$$\begin{bmatrix} v(1) \\ v(2) \\ \vdots \\ v(n) \end{bmatrix}$$

It can be solved directly,

$$v = R + \gamma P v$$

$$(1 - \gamma P) v = R$$

$$v = (1 - \gamma P)^{-1} R$$

$O(n^3)$  for  $n$  states

To solve large MDPs

- DP

- Monte-Carlo

- Temporal-Difference Learning

MDP is a Markov Reward Process with decisions

$$(S, A, P, R, \gamma)$$

$$P_{ss'}^a = P[S_{t+1} = s' | S_t = s, A_t = a]$$

$$R_s^a = \mathbb{E}[R_{t+1} | S_t = s, A_t = a]$$

Policy  $\pi$  is a distribution over actions given states.

$$\pi(a|s) = P[A_t = a | S_t = s]$$

MDP policies depend on current state

state value

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

action value

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

Value function

$$V_\pi = E_\pi [R_{t+1} + \gamma v_\pi(s_{t+1}) | S_t = s]$$

Q function

$$q_\pi = E_\pi [R_{t+1} + \gamma q_\pi(s_{t+1}, a_{t+1}) | S_t = s, A_t = a]$$

Optimal value function

$$V^*(s) = \max_\pi v_\pi(s)$$

$$q^*(s, a) = \max_\pi q_\pi(s, a)$$

optimal policy

$\pi \geq \pi'$  if  $V_\pi(s) \geq V_{\pi'}(s)$  vs.

Theorem

- For any MDP, there is an optimal policy better than or equal to all the policy.

$$\pi^* > \pi \text{ vs}$$

- All optimal policies achieve the optimal value function's action-value fn.

$$V_{\pi^*}(s) = v^*(s)$$

$$V_{\pi^*}(s, a) = q^*(s, a)$$

An optimal policy can be found by maximum over  $q^*(s, a)$

$$\pi^*(a|s) = \begin{cases} 1 & \text{if } a = \operatorname{argmax}_a q^*(s|a) \\ 0 & \text{otherwise} \end{cases}$$

Bellman Optimality Eq for  $v^*$

$$v^*(s) = \max_a q^*(s|a)$$

$$q^*(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v^*(s')$$

$$\therefore v^*(s) = \max_a R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a v^*(s')$$

$$q^*(s, a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a \max_a q^*(s'|a)$$

These are non linear.

To solve it

- Value iteration
  - policy iteration
  - Q-learning
- } DP-based

Finally - Sarsa - ~~Value function~~ ~~and value function~~ ~~it's~~ ~~function~~ ~~function~~

## Planning by DP

- Sequential a temporal aspect to the problem.
- Programming optimising a 'program' - a policy.
- Breaking down into subproblems and solving them.
- Can be applied on problems which have properties
  - optimal substructure  
Can solve an overall problem by solving optimal solutions of its different pieces.
  - Overlapping subproblem  
Subproblems occur again & again

- MDP satisfies both properties
- Value function stores and retrieves solutions
- For prediction :
  - i/p:  $MDP(S, A, P, R, \gamma)$ ,  $s_i \in \pi$
  - $MRP(s, P^\pi, R^\pi, \gamma)$
- o/p - value function  $v_\pi$
- For control :
  - i/p -  $MDP(S, A, P, R, \gamma)$
  - optimal  $v_\pi, \pi^*$

### Policy Evaluation

Problem - To calculate a given policy  $\pi$

Soln

- iterative application of Bellman Expectation Eq.

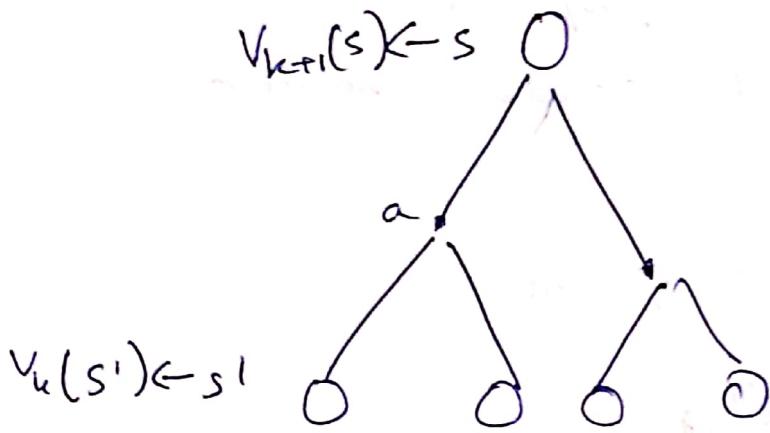
$$v_1 \rightarrow v_2 \rightarrow \dots \rightarrow v_\pi$$

- Using synchronous backups

. At every iteration  $k+1$

. For all states  $s \in S$

. Update  $v_{k+1}(s)$  from  $v_k(s')$

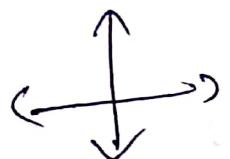


~~Repeat~~ place the current value function at the every leaf  $s$

compute new  $V_k$  using Bellman eq (Multiply by Probability & sum it up)  
and repeat this

This process will converge on the true value fn.

on a small grid



action

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

$r_{2-1}$   
for all.

1 terminal state (2 squares)

find the avg steps taken to reach terminal with uniform random policy.

$$\pi(n_1 \cdot) = \pi(e_1 \cdot) = \pi(s_1 \cdot) = \pi(w_1 \cdot) = 0.25$$

$v_k$  for the random policy

$k=0$

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

(our estimate of every single state  
on how much  
steps to reach  $\times$ )  
(start with observation)

Then for each state  
compute Bellan expectation

$k=1$

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

at every state, any step will give only -1 reward

$k=2$

$\sigma$	-1.75	2	-2
-2	-2	-2	-2
-2	-7	-7	-1.75
-2	-2	-1.75	0

$\hat{f} \leftarrow \underline{\sigma}$ .

$$\begin{aligned}
 &= (-1 + -1) + (1 + -1) + (-1 + -1) \\
 &\quad + \underline{(1 + 0)} \quad -1.75 \\
 &\quad \quad \quad 4
 \end{aligned}$$

$k=\infty$       0      -14      -20      -22

-14      -18      -20      -20

-20      -20      -18      -14

-22      -20      -14      0

greedy policy wrt  $v_k$

look at  $v$  and predict what value to take.

In 3 iteration it can give optimal policy.

So value fn helps us figure out policy.

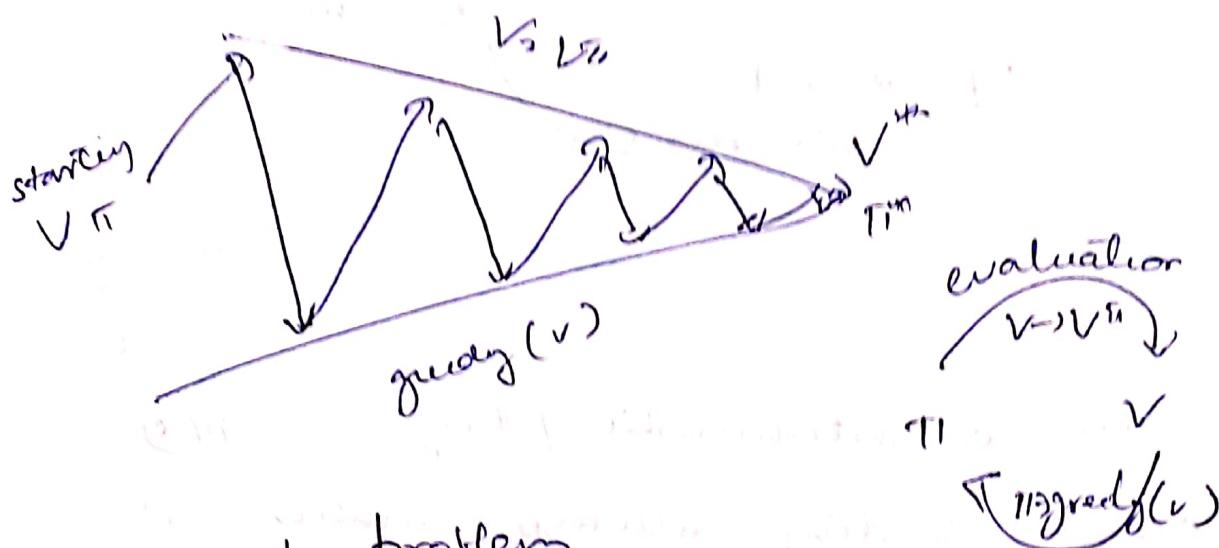
## policy iteration

- To find the best policy
- Evaluate the policy

$$v_{\pi}(s) = \mathbb{E}[R_{t+1} + \gamma R_{t+2} + \dots | s_t = s]$$

- Improve the policy by acting greedily

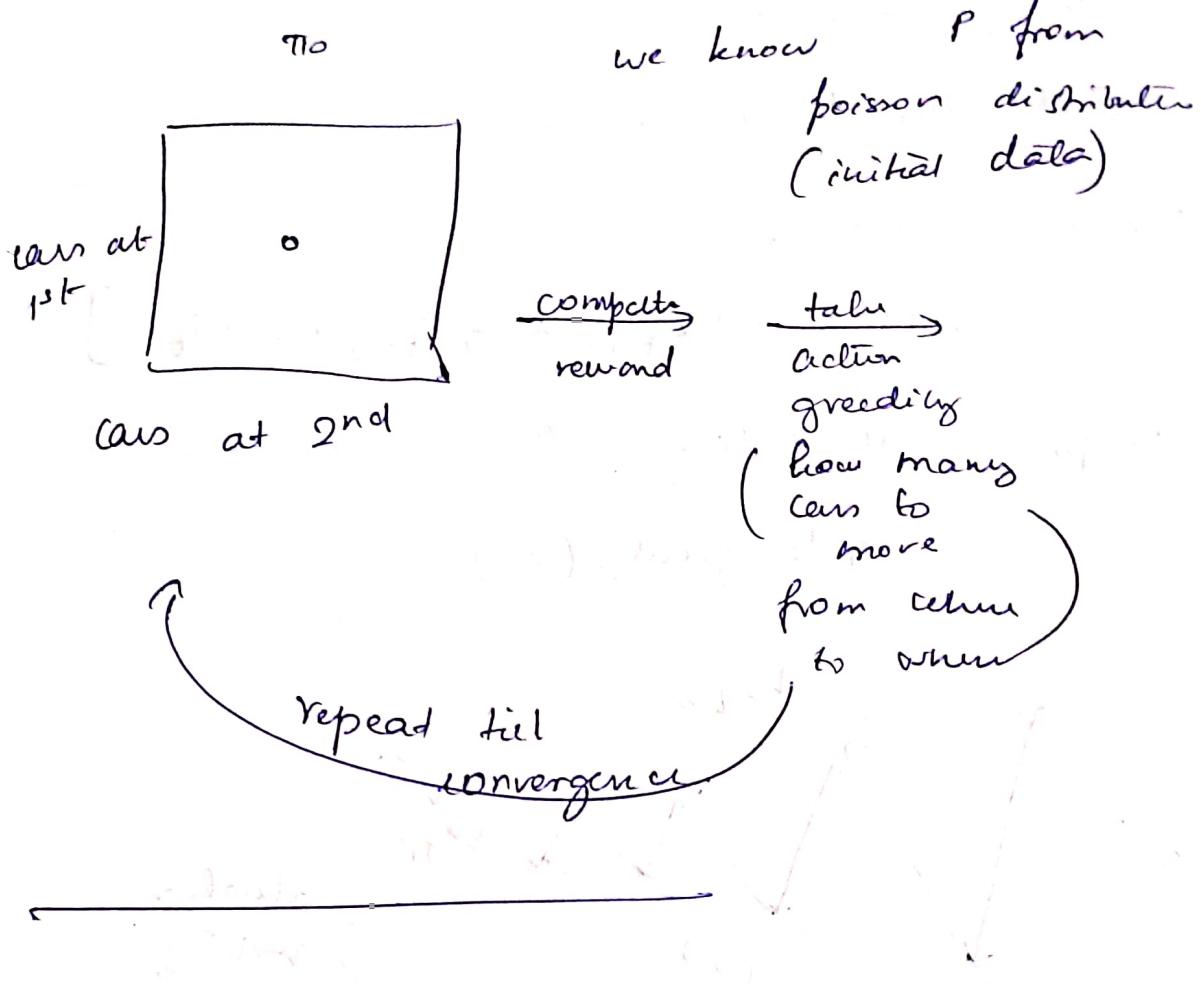
$$\pi^{\text{greedy}}(v_{\pi})$$



## Car rental problem

States

- 2 locations, max of 20 cars
- Actions - Move up to 5 cars b/w locations
- Reward - \$10 for each rental overnight
- Transition - Cars returned & required randomly
  - 1st location  $\rightarrow$  avg reqn, avg return
  - 2nd location  $\rightarrow$  higher with  $\geq$  returns



- Consider a deterministic policy  $\pi = \pi(s)$   
(acting greedily always makes it deterministic)
- This improves value over one step:

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

$$\begin{aligned} V_{\pi}(s) &\leq V_{\pi'}(s) \geq E_{\pi'}[R_{t+1} + \gamma V_{\pi}(s_{t+1}) | s_t] \\ &\leq E_{\pi'}[R_{t+1} + \gamma V_{\pi'}(s_{t+1}, \pi'(s_{t+1})) | s_t] \end{aligned}$$

$$\leq \mathbb{E}_{\pi^1} [R_{t+1} + \gamma R_{t+2} + \dots | S_t = s] = v_{\pi^1}(s)$$

This will make ~~most~~<sup>it</sup> step attract when  
At some point it might stop.  
if improvements stop,

$$q_{\pi^1}(s, \pi'(s)) > \max_a q_{\pi^1}(s, a) = q_{\pi^1}(s, \pi(s)) = v_{\pi^1}(s)$$

↓

This satisfies BOE

$$v_{\pi^1}(s) = \max_{a \in A} q_{\pi^1}(s, a)$$

Policy iterations solves MDPs.

### Modified policy iteration

- Should we include a stopping condition?
- Stop after k iteration?
- Why not update policy every iteration?

### Value iteration

Principle of optimality:

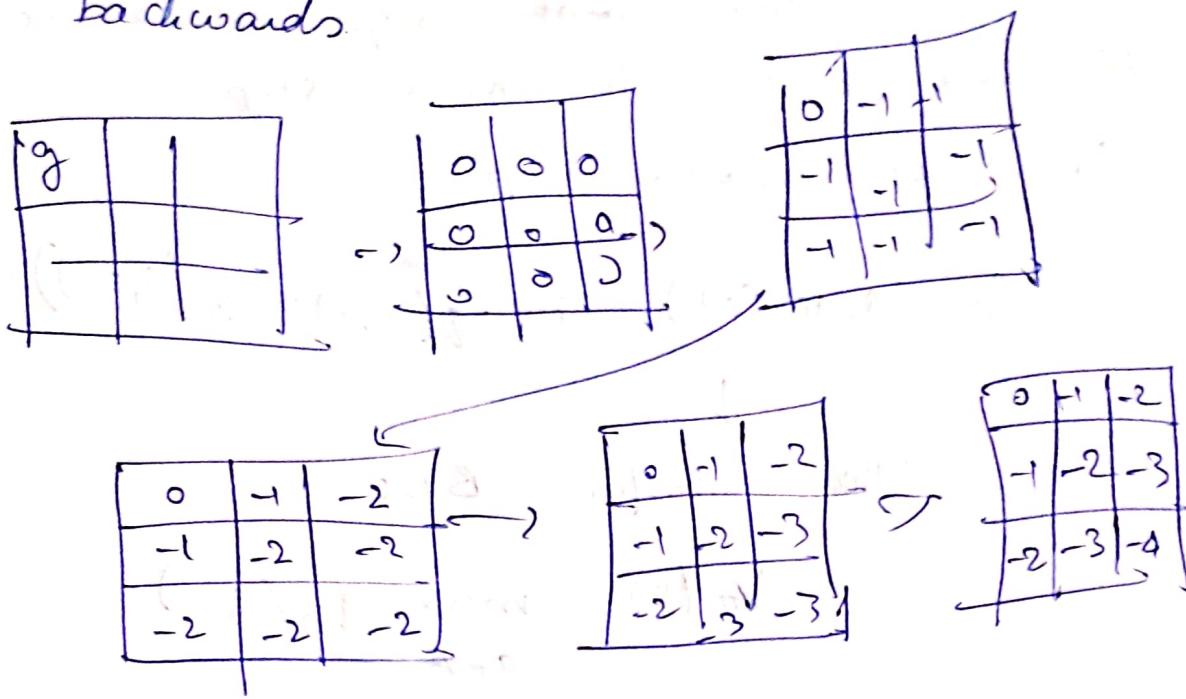
A policy  $\pi(a|s)$  achieves the optimal value from state  $s$   $v(s) = \pi_x(s)$  if and only if

- For any state  $s'$  reachable from  $s$
- $\pi$  achieves optimal value from state  $s'$

Intuition:

Start with final rewards  
backwards

& work



one step lookahead from the  
first terminal.

value iteration

- to find optimal policy  $\pi$

Asynchronous DP

- backs up states individually
- will converge if all states are selected at some point.
- in place DP
- prioritized sweeping
- Real time DP

- In-plan → forget about old & new value func
- just use the latent values,

### Prioritized Sweeping

- priority queue to see which states are better than others.
- magnitude of error in Bellman equation b/w states

- Real-time - Select States that agent selects

DP uses full-width backups,

(computation on all nodes)

- so sample backups

<u>Model</u>	<u>free</u>	<u>Volatil</u>	<u>Prediction</u>
--------------	-------------	----------------	-------------------

### Monte Carlo learning

- look at complete episodic memory.
- learn  $v_{\pi}$  from episodes of experience from

$$b_t = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$$

$$v_{\pi}(s) = E_{\pi}[b_t | s_t = s]$$

In MC, empirical mean is taken instead of expectation.

### First Visit MC policy evaluation

- consider the first time one visits a state
- increment a counter ( $N(s)$ )
- add up total return ( $S(s) = s(s) + b_t$ )
- take mean from that time
- By law of large numbers,  $V(s) \rightarrow v_{\pi}(s)$

as  $N(s) \rightarrow \infty$

### Every Visit MC policy Evaluation

- Consider every visit to the state & sum it up

## Incremental mean

$$\mu_k = \frac{1}{k} \sum_1^k x_j$$

$$= \frac{1}{k} \left[ x_k + \sum_{j=1}^{k-1} x_j \right]$$

$$= \frac{1}{k} [x_k + (k-1)\mu_{k-1}]$$

$$\mu_{k-1} + \frac{1}{k} (x_k - \mu_{k-1})$$

↙      ↓      ↓  
 total    new    previous step mean  
 steps    step    value

## Incremental MC updates

For each state  $s_t$  with  $g_t$

$$N(s) = N(s_t) + 1$$

$$V(s_t) = V(s_t) + \frac{1}{N(s_t)} (g_t - V(s_t))$$

To forget old episodes replace  $\frac{1}{N(s_t)}$  by a constant  $\alpha$ .

## Temporal-Difference Learning

- learns from incomplete episodes.
- It updates a guess towards a guess

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

- Simplest TD algorithm, TD(0)

update value  $V(s_t)$  towards  $\alpha$  estimated return  $R_{t+1} + \gamma V(s_{t+1})$

$\downarrow$                              $\downarrow$

immediate reward      discounted value of next step.

$$V(s_t) = V(s_t) + \alpha (R_{t+1} + \gamma V(s_{t+1}) - V(s_t))$$

$\underbrace{\qquad\qquad}_{\text{TD target}}$

$\underbrace{\qquad\qquad}_{\text{TD error}}$

### TD vs MC

- TD can learn before final outcome
  - MC should wait till end
- 
- TD can learn when it may not see final outcome
  - MC can learn only from complete sequences
- TD works in non-terminating env
  - MC works in terminating env

## Bias, Variance tradeoff

- $G_t$  is unbiased of  $v_{\pi}(s_t)$  in MC.
- True TD is unbiased estimate of  $v_{\pi}(s_t)$   
(if  $v_{\pi}$  is known)
- TD target with  $V(s_{t+1})$  is biased estimate  
of  $v_{\pi}(s_t)$
- TD target has less variance than return.  
(return depends on lots of terms)  
(TD target depends on 1 random transition)
- MC has good convergence due to no bias
- Not sensitive to initial values
- TD more sensitive to initial values
- TD much more efficient than MC.
- MC converges to solution with min MSE
- TD converges to soln of max likelihood Markov model
- MC more effective in the non Markov case
- TD exploits Markov property, much better

Bootstrapping - update involves

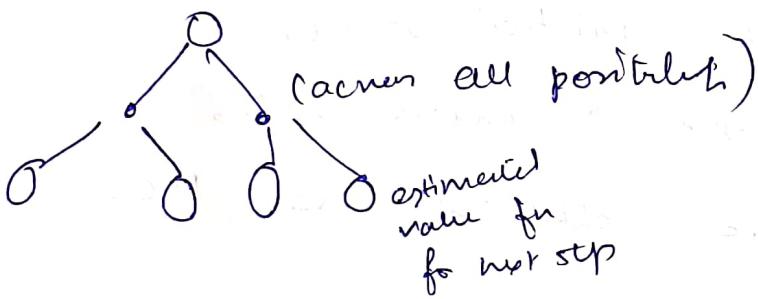
- MC does not bootstrap
- TD, DP bootstraps

estimate  
(estimated  
Value for  
of next step)

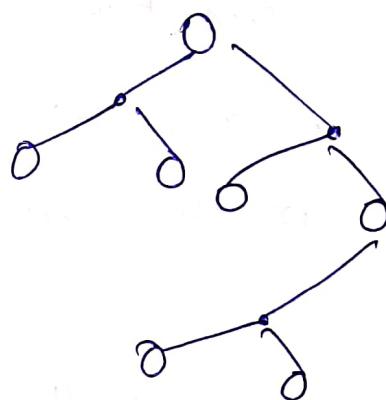
### Sampling

- MC samples
- DP does not (every possibility considered)
- TD samples

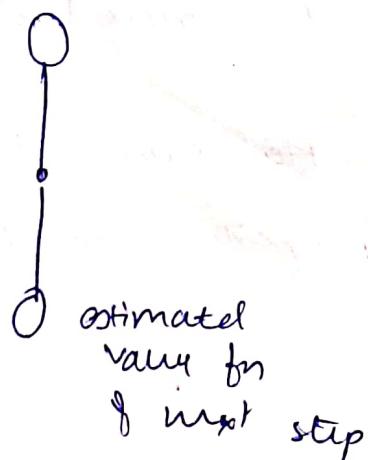
### DP



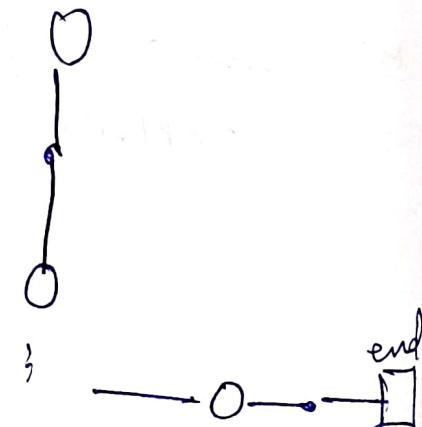
### Exhaustive search



### TD



### MC



## TD( $\lambda$ )

n steps of look ahead and estimating value.

n = total steps, TD = MC

$$G_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n}$$

$$V(s_t) = V(s_t) + \alpha \cdot (G_t^{(n)} - V(s_t))$$

To select best n:

average over n-step returns over diff n

Eg: avg of 2-step & 4-step

$$\frac{1}{2} G_t^{(2)} + \frac{1}{2} G_t^{(4)}$$

$\lambda$ -Return

geometrically weighted avg of all n.

$$G_t^\lambda = (1-\lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_t^{(n)}$$

geometric weights seem efficient

(less memory)

This is forward view

(computation same as TD(0))

### Backward view

- keeps eligibility trace for every state  $s$

↑  
to error

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

$$V_t(s) = V(s) + \alpha \delta_t E_t(s)$$

when  $\gamma = 0$ ,

it is TD(0)

when  $\gamma = 1$ ,

credit is diffused until end of episode.

### Model Free Control

- On policy & Off-policy.

#### On-policy MC control

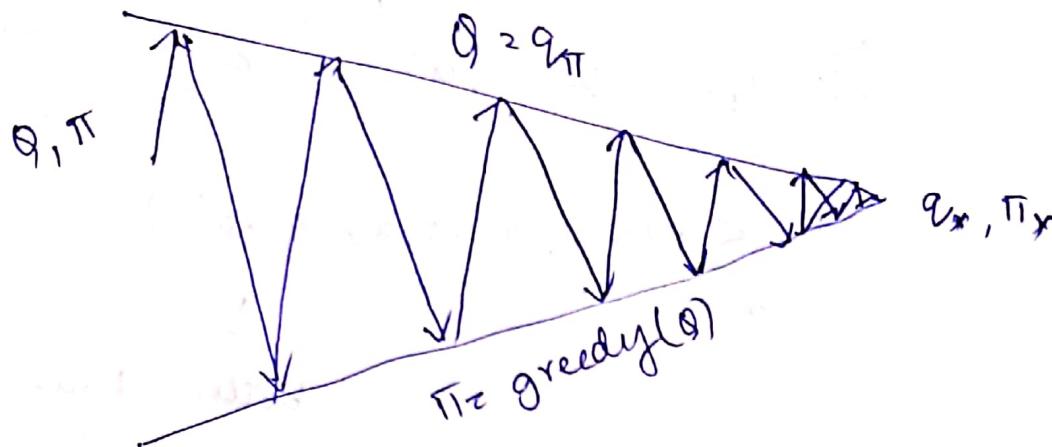
will replacing DP with MC work?

No, because

- no exploration (acting greedily will ignore a lot of state spaces)
- greedy policy over  $V(s)$  act greedily. (no transmission requires model of MDP to matrix available)

Alternative is to use action value functions  
 $Q(s, a)$  is model free

$$\pi'(s) = \operatorname{argmax}_{a \in A} Q(s, a)$$



### $\epsilon$ -greedy exploration

- simplest idea

- All  $m$  actions are tried with non-zero prob
- with prob  $1-\epsilon$ , choose the greedy action
- with prob  $\epsilon$ , choose an action at random
- This guarantees to continue  $\epsilon$  explore everything
- with  $\epsilon$ -greedy, there will be step of policy improvement.

$\epsilon$ -greedy Policy Improvement proof.

that we will take each action

$$q_{\pi'}(s, \pi'(s)) = \sum_{a \in A} \pi'(a|s) q_{\pi}(s|a)$$

$\downarrow$

one step by new policy

$$= \epsilon \underset{\text{prob of taking every action}}{\underset{\text{any action}}{\sum}} q_{\pi}(s, a) + (1-\epsilon) \max_{a \in A} q_{\pi}(s, a)$$

↑ prob of taking greedy act.

$$\geq \epsilon \underset{\text{any action}}{\underset{\text{any action}}{\sum}} q_{\pi}(s, a) + (1-\epsilon) \sum \frac{\pi(a|s) - \epsilon/m}{1-\epsilon} q_{\pi}(s, a)$$

$$\sum_{a \in A} \pi(a|s) q_{\pi}(s, a) = v_{\pi}(s)$$

$\downarrow$

better than value of original policy

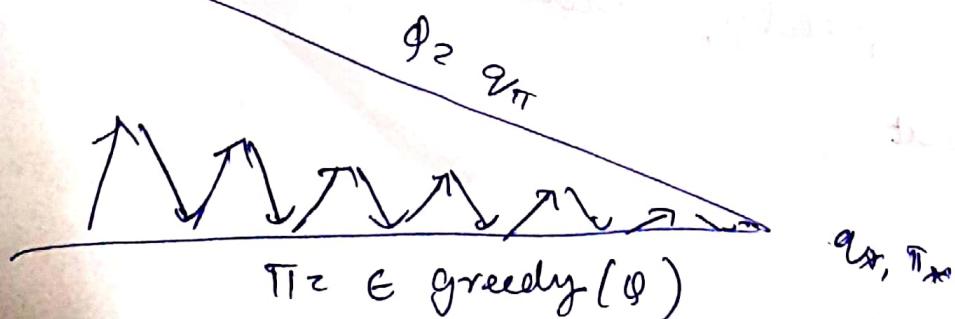
this sum max is considered

of all the actions greater than weighted

$$\therefore v_{\pi'}(s) \geq v_{\pi}(s)$$

update every single episode,

- act greedily wrt fresh estimate



## greedy in limit with Infinite Exploration (GLIE)

- All action states are infinitely explored many times.
- Policy converges on being greedy.

$\epsilon$ -greedy is GLIE if  $\epsilon$  reduces to zero

$$\text{at } \epsilon_k = \frac{1}{k}$$

## GLIE MC Control

- Sample  $k^{\text{th}}$  episode  $\{s_1, a_1, r_2, \dots, s_T\} \sim \pi$
- For each state & action

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

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \frac{1}{N(s_t, a_t)} (r_t + \gamma Q(s_{t+1}, a_{t+1}))$$

- Improve policy based on new action values

$$\epsilon = 1/k$$

$$\pi \leftarrow \epsilon\text{-greedy}(Q)$$

TD control advantages over MC.

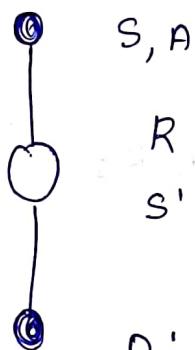
- lower variance
- Online
- Incomplete sequences

Natural idea:

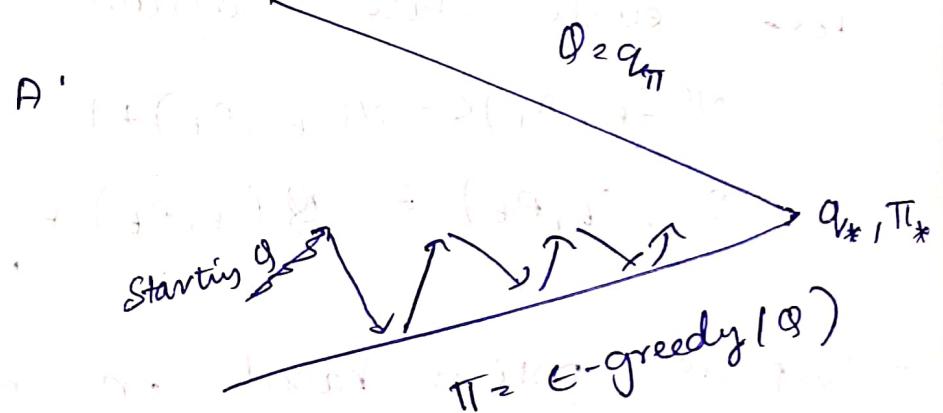
use TD instead of MC in control loop.

- Apply TD to  $Q(s, a)$
- Use  $\epsilon$ -greedy
- Update every time step.

### Updating Action Value fns with Sarsa



$$Q(s, a) = Q(s, a) + \alpha (R + \gamma Q(s', a') - Q(s, a))$$



### Pseudo code

- Initialize  $Q(s, a)$  arbitrarily
- Repeat (for every episode)
  - Initialize  $s$
  - Choose  $A$  from  $s$  using policy driven for  $Q$
  - Repeat (for each step)
    - take action, observe  $R, s'$
    - Choose  $A'$  from  $s'$  using  $\epsilon$ -greedy
    - update  $Q$
- terminal? update  $s, a$

- bLIE Sarsa will converge.

- Robbins - Monro sequence of step-sizes  $\alpha_t$   
 $\sum \alpha_t = \infty$   
 $\sum \alpha_t^2 < \infty$

n-step sarsa

Same as n-step TD,

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha (q_t^{(n)} - Q(s_t, a_t))$$

Sarsa  $\Rightarrow$  (same as TD) (this is not online)

$q^\lambda$  returns combines all n-step  $\pi$ . Q-return

$$q_t^{(n)}$$

Using weight  $(1-\lambda)^{\lambda^{n-1}}$

$$q_t^\lambda = (1-\lambda) \sum_{i=1}^{\infty} \lambda^{n-i} q_t^{(n)}$$

- Forward view Sarsa( $\lambda$ )

$$Q(s_t, a_t) = Q(s_t, a_t) + \alpha (q_t^\lambda - Q(s_t, a_t))$$

- Backward view Sarsa( $\lambda$ )

- same as TD( $\lambda$ ) Backward

- But Sarsa( $\lambda$ ) has one ET for each state action.

$$E_0(s, a) = 0$$

$$E_t(s, a) = \gamma \lambda E_{t-1}(s, a) + I(s_t=s, a_t=a)$$

## Off-policy learning.

- Learning from others.
- Learning from a behavior policy  $\mu$
- Re-use experience from old policies
- Learn about optimal policy while following exploratory policy
- Learn about multiple policies while following one policy.

## Important Sampling.

- Estimate expectation of a diff distribution

$$\begin{aligned} E_{x \sim p}[f(x)] &= \sum p(x)f(x) \\ &= \sum q(x) \frac{p(x)}{q(x)} f(x) \end{aligned}$$

$E_{x \sim q}\left[\frac{p(x)}{q(x)} f(x)\right]$

## Important Sampling from MCMC

- have to multiply IS ratio to every step
- all returns generated from  $\mu$  to evaluate  $\pi$
- extremely high variance (policies don't match over sequences)
- Not used in practice.

## Importance Sampling for Off-policy TD

- Multiply for one step
- Lower Variance
- Policies need to be similar only for one step.

## Q-learning

- Specific to TD(0) or SARSA(0)
- use action states  $Q(s, a)$
- select next action according to behavior policy
- Consider alternative success action if we followed target policy.  
 $A_{t+1} \sim \pi(\cdot | s_t)$   
 $A' \sim \pi(\cdot | s_t)$
- And update  $Q(s_t, A_t)$  towards value of alternative action

$$Q(s_t, A_t) = Q(s_t, A_t) + \alpha (R_{t+1} + \gamma Q(s_{t+1}, A') \cdot Q(s_t, A))$$

$\downarrow$                                      $\downarrow$   
μ policy                              target policy

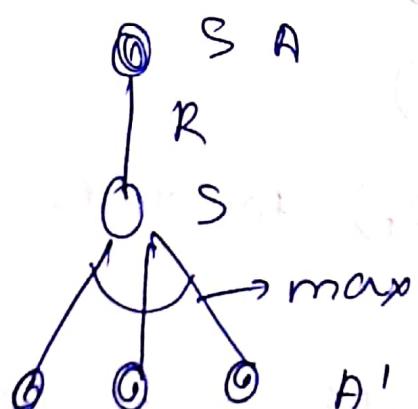
- No Importance policy.

- target policy is a greedy policy
- both  $\mu$  &  $\pi$  can improve
- $\pi$  is made greedy wrt  $Q$ .  

$$\pi(s_{t+1}) = \arg\max Q(s_{t+1}, a')$$
- $\mu$  is  $\epsilon$ -greedy wrt  $Q$
- 
- $$q_t = R_{t+1} + \gamma Q(s_{t+1}, a')$$
- $$= R_{t+1} + \max_{a'} Q(s_{t+1}, a')$$

estimated at time  $t$ :  $R_{t+1} + \max_{a'} Q(s_{t+1}, a')$

and update  $Q$   $\downarrow$   
 updates till bit on max  $Q$   
 value it can take



## Value function Approx

### Large Scale RL

- Back gammon  $10^{20}$  states
- go  $10^{170}$  states
- helicopter continuous state space.

so far

- lookup table for  $Q(s, a)$  or  $v(s)$  for every  $s$

### Problem with large MDP

- too many states & actions to store in memory
- too slow to learn everything individually.

### Solution

- Value fn approx

$$\hat{v}(s, w) \approx v_{\pi}(s)$$

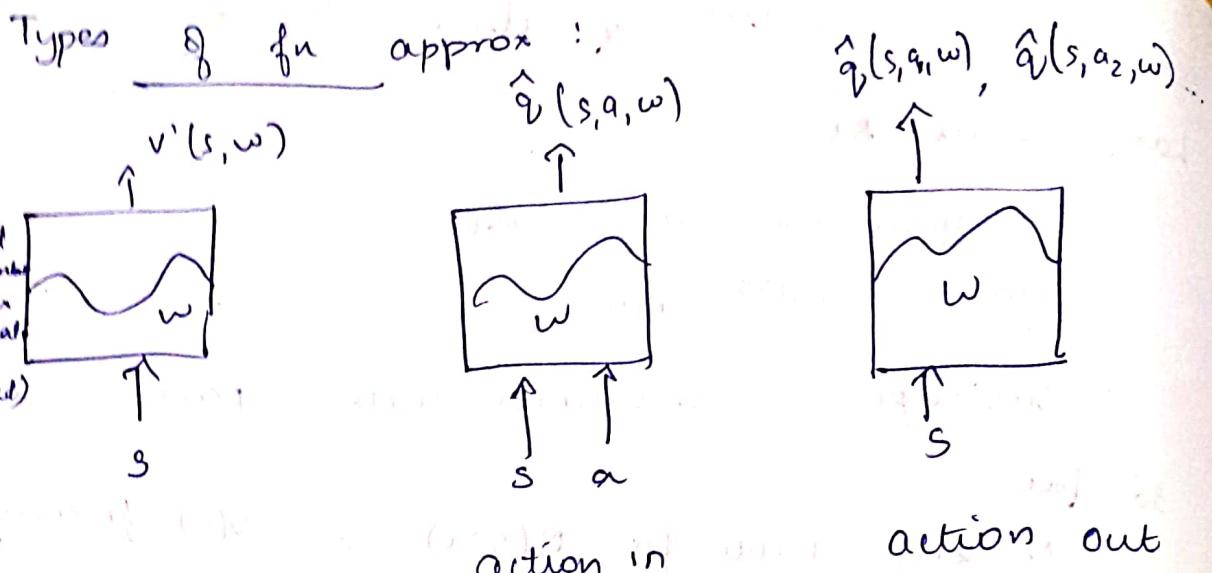
↓

estimate of  $v_{\pi}(s)$

$w \rightarrow$  weights

$$\hat{q}(s, a, w) \approx q_{\pi}(s, a)$$

- generalizing seen states to unseen states
- update parameters  $w$  using TD or MC



### function approx.

- Linear combinations of features
- Neural net w/ w as weight of } differential
- decision tree
- nearest neighbour
- Fourier / Wavelet basis

### Incremental methods

gradient descent:

$\delta(\omega)$  diff param of  $\omega$

$$\nabla_{\omega} J(\omega) = \begin{pmatrix} \frac{\partial J(\omega)}{\partial \omega_1} \\ \vdots \\ \frac{\partial J(\omega)}{\partial \omega_n} \end{pmatrix}$$

To find local minima of  $J(\omega)$

Adjust param  $\omega$  in direction of -ve gradient

$$\Delta \omega = -\frac{1}{2} \alpha \nabla_{\omega} J(\omega)$$

Value for approx using SGD

- find  $w$  minimising MSE  $J(w)$   $\hat{v} \in V_\pi$

$$J(w) = E_\pi [(V_\pi(s) - \hat{v}(s, w))^2]$$

$$\therefore \Delta w = \alpha E_\pi [(V_\pi(s) - \hat{v}(s, w)) \nabla_w \hat{v}(s, w)]$$

SGD samples the gradient

$$\Delta w = \alpha (V_\pi(s) - \hat{v}(s, w)) \nabla_w \hat{v}(s, w)$$

Feature vector

Represent state by FV

$$s = \begin{pmatrix} x_1(s) \\ \vdots \\ x_n(s) \end{pmatrix}$$

- e.g.: distance from landmarks
- : trends in stock market
- : piece config in chess

Linear VFA

- linear combination of features

$$\hat{v}(s, w) = x(s)^T w = \sum_j x_j(s) w_j$$

$$J(w) = E_\pi [(V_\pi(s) - x(s)^T w)^2]$$

$$\nabla_w \hat{v}(s; w) = x(s)$$

$$\Delta w = \alpha (V_\pi(s) - \hat{v}(s, w)) x(s)$$

$$\text{update} = (\text{step-size}) \times (\text{error}) \times (\text{feature value})$$

But there are no supervising, only rewards.

In practise, substitute a target for  $v_T(s)$

For MC, target is  $g_t$

$$\therefore \Delta w = \alpha(g_t - \hat{v}(s_t, w)) \nabla_w \hat{v}(s_t, w)$$

For TD, target is  $R_{t+1} + \gamma \hat{v}(s_{t+1}, w)$

$$\Delta w = \alpha(R_{t+1} + \gamma \hat{v}(s_{t+1}, w) - \hat{v}(s_t, w)) \nabla_w \hat{v}(s_t, w)$$

For TD( $\lambda$ ), target is  $\lambda$ -return.  $g_t^\lambda$ .

$$\Delta w = \alpha(g_t^\lambda - \hat{v}(s_t, w)) \nabla_w \hat{v}(s_t, w)$$

### MC with VFA

- return is used to build some training data

$$(s_1, g_1), (s_2, g_2), \dots (s_T, g_T)$$

- converges to global optimum

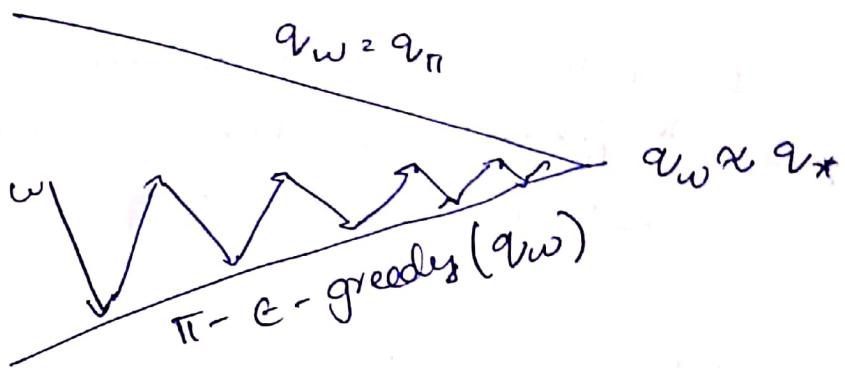
### TD with VFA

target is a sample of true value  $v_T(s_t)$

$$(s_1, R_2 + \gamma \hat{v}(s_2, w)), (s_2, R_3 + \gamma \hat{v}(s_3, w)), \dots (s_{T-1}, R_T)$$

- converges close to global optimum.

Control with VFA



Action value fn Approx

$$\hat{q}_v(s, a, \omega) = a_{\pi}(s, a)$$

$$J(\omega) \rightarrow E_{\pi} [(q_{\pi}(s, a) - \hat{q}_v(s, a, \omega))^2]$$

$$x(s, a) \rightarrow \begin{pmatrix} x_1(s, a) \\ \vdots \\ x_n(s, a) \end{pmatrix}$$

Same as before for everything

Batch RL

- Grid is not sample efficient
- Experience D consisting of

$$D = \{(s, v_i^{\pi}), (s_2, v_2^{\pi}) \dots\}$$

Sample  $s, v$  from D  
apply  $\rightarrow$  grid

(state, value) pairs

Exp Replay is solved via least square prediction

## DQN

- uses exp replay & fixed Q-targets
- for which move towards q value in old  $w$  (keeping it fixed) & then make  $w$  mean equal to new  $w$ .

## Policy gradient

directly parameterise policy instead of  $v_\pi$  or  $q_\pi$

---

Advantages of policy based or value based

- Better convergence
- Effective in high dim
- Can learn stochastic policies

## Pis Ad

- Typically converge to lower than upper
- Evaluating a policy is typically inefficient & high variance

## Policy obj functions

measure of quality of policy  $\pi_\theta$ ?

- In episode env, we start value
- $$\pi_1(0) = v_{\pi_0}(s_1) = E_{\pi_0}(v_1)$$

- continuing on, we avg value
- $$J_{av}(\theta) = \sum_s d^{\pi^\theta}(s) V^{\pi^\theta}(s)$$

- avg reward per time step

$$J_{avr}(\theta) = \sum_s d^{\pi^\theta}(s) \sum \pi_\theta(s, a) R_s^a$$

To optimise time objectives:

- gradient free
  - Hill climbing
  - genetic algo
  - ~~- simplex / Amoeba / Nelder mead~~
- gradient based
  - GD
  - conjugate grad
  - Quasi-Newton.

### policy grad

- $J(\theta)$  - obj fn
- gradient ascent to maximise local max  $\nabla_\theta J(\theta)$

$$\Delta\theta = \alpha \nabla_\theta J(\theta)$$

Policy grad

$$\nabla_\theta J(\theta) = \begin{pmatrix} \frac{\partial J(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_n} \end{pmatrix}$$

## Computing grads by finite diff

$$\frac{\partial J(\theta)}{\partial \theta_k} \approx \frac{J(\theta + \epsilon u_k) - J(\theta)}{\epsilon}$$

- simple, noisy, inefficient
- works even if it is not diff.

## Analytically compute grads

- score fn.
- assume  $\pi_\theta$  is diff when picking action
- we know  $\nabla_\theta \pi_\theta(s, a)$
- Likelihood ratio :

$$\nabla_\theta \pi_\theta(s, a) = \nabla_\theta \log \pi_\theta(s, a)$$

Score fn is  $\nabla_\theta \log \pi_\theta(s, a)$

## Softmax policy

$$\pi_\theta(s, a) \propto e^{\phi(s, a)^T \theta}$$

$$\text{score fn} = \phi(s, a) - E_{\pi_\theta} [\phi(s, \cdot)]$$

$\hookrightarrow$  feature we took       $\hookleftarrow$  avg of all features

gaussian policy:

(on cont. action space)

parameters  $\mu$  of the gaussian

$$\mu(s) = \phi(s)^T \theta$$

Variance can be fixed or parameterized

$$\sigma \sim N(\mu(s), \sigma^2)$$

Score function =

$$\nabla_{\theta} \log \pi_{\theta}(s, a) = \frac{(a - \mu(s)) \phi(s)}{\sigma^2}$$

for one-step MDP

- start at state  $s$  and  $d(s)$
- terminate after one step, reward  $r \sim P_{s,a}$
- Using likelihood ratios to compute the policy gradient

$$J(\theta) = E_{\pi_{\theta}}[r]$$

$$= \sum d(s) \sum \pi_{\theta}(s, a) R_{s,a}$$

$$\nabla_{\theta} J(\theta) = \sum d(s) \sum \pi_{\theta}(s, a) \nabla_{\theta} \log \pi_{\theta}(s, a) R_{s,a}$$

$$= E_{\pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, \theta) r]$$

for multi step

- replace  $r$  with long term value  $\mathbb{Q}^{\pi}(s, a)$

$$J \approx J_1, J_{\text{max}} \text{ or } \frac{1}{1-\gamma} J_{\text{avg}}$$

$$\text{all } \rightarrow E[\nabla_{\theta} \log \pi_{\theta}(s, \theta) \mathbb{Q}^{\pi_{\theta}}(s, \theta)]$$

## Monte-Carlo policy grad (Reinforce)

- update params by SGD
- use Pg to sample exp.

$$\Delta \theta_t = \alpha \nabla_{\theta} \log \pi_{\theta}(s_t, a_t) V_t$$

for each episode  $\{s_1, a_1, r_2, \dots\} \sim \pi^{\theta}$  do

$$\theta_t \leftarrow \theta + \Delta \theta_t$$

## Reducing Var using a critic

- MC Pg has high var
- use critic to estimate  $Q$

$$Q_w(s, a) \approx Q^{\pi^{\theta}}(s, a)$$

Action critic  $\rightarrow$  2 sets of params

Critic  $\rightarrow$  for action value fn, params

actor  $\rightarrow$  for policy, params  $\theta$ , in direction suggested by critic

## Approx policy grad combining both:

$$\nabla_{\theta} J(\theta) \approx E_{\pi^{\theta}} [\nabla_{\theta} \log \pi_{\theta}(s, a) Q_w(s, a)]$$

$$\Delta \theta = \alpha \nabla_{\theta} \log \pi_{\theta}(s, a) Q_w(s, a)$$

To estimate  $Q_w$ :

- policy evaluation: how good is  $\pi^{\theta}$  under  $\theta$ ?
- MC or TD or TD( $\lambda$ )

## Bias in Actor-critic

- Approx. pg intro bias

### Reducing Var using Baseline

- Subtract a Baseline for  $B(s)$  from policy  $g$ .
- This can change var w/o changing exp.

$$E_{\pi_0} [\nabla \log \pi_0(s, a) B(s)]$$

$$= \sum_a d^{\pi_0}(s) \sum_a \nabla_\theta \pi_0(s, a) B(s)$$

$$\approx \sum_a d^{\pi_0} B(s) \nabla_\theta \sum_a \pi_0(s, a)$$

gradient becomes  
 $\leftarrow$  (prob dist. adds to 1)  $\rightarrow = 0$

we can use state value fn as  $B(s)$

$$B(s) = v^{\pi_0}(s)$$

rewriting pg using  $\alpha$  the advantage  
fn  $A^{\pi_0}(s, a)$

$$A^{\pi_0}(s, a) = q^{\pi_0}(s, a) - v^{\pi_0}(s)$$

$$\nabla_\theta J(\theta) = E_{\pi_0} [\nabla \log \pi_0(s, a) A^{\pi_0}(s, a)]$$

depending on whether it is tve or -ve,  
it changes the policy.

## Estimating the advantage fn:

- Critic learns  $\hat{Q}$  and  $V$ .

$$V_v(s) \approx V^{\pi_0}(s)$$

$$Q_w(s, a) = Q^{\pi_0}(s, a)$$

$$A(s, a) = Q_w(s, a) - V_v(s)$$

update both using TD learning

- Commonly used variant

- TD error considered as a sample of the adv fn.

$$\delta^{\pi_0} = r + \gamma V^{\pi_0}(s') - V^{\pi_0}(s)$$

$$E_{\pi_0}[\delta^{\pi_0}|s, a] \approx A^{\pi_0}(s, a)$$

$$\therefore \nabla_\theta J(\theta) \approx E_{\pi_0}[\nabla_\theta \log \pi_0(s, a) \delta^{\pi_0}]$$

- can use approx TD error

$$\delta_V = r + \gamma V_v(s') - V_v(s)$$

- only 1 set of critic params need

Critics at diff time scales.

For MC

$$\Delta \theta \approx \alpha(V_t - V_\theta(s)) \phi(s)$$

For TD(0)

$$\Delta \theta \approx \alpha(r + \gamma V(s') - V_\theta(s)) \phi(s)$$

For Forward TD(λ)

$$\Delta \theta \approx \alpha(V_t^\lambda - V_\theta(s)) \phi(s)$$

For backward TD(λ)

$$\delta_t = r_{t+1} + \gamma V(s_{t+1}) - V(s_t)$$

$$e_t = \lambda e_{t-1} + \phi(s_t)$$

$$\Delta \theta = \alpha \delta_t e_t$$

saw for actors

Alternative Policy gradient directions

- compatible for Approx.
- features for critic can scores of policy.

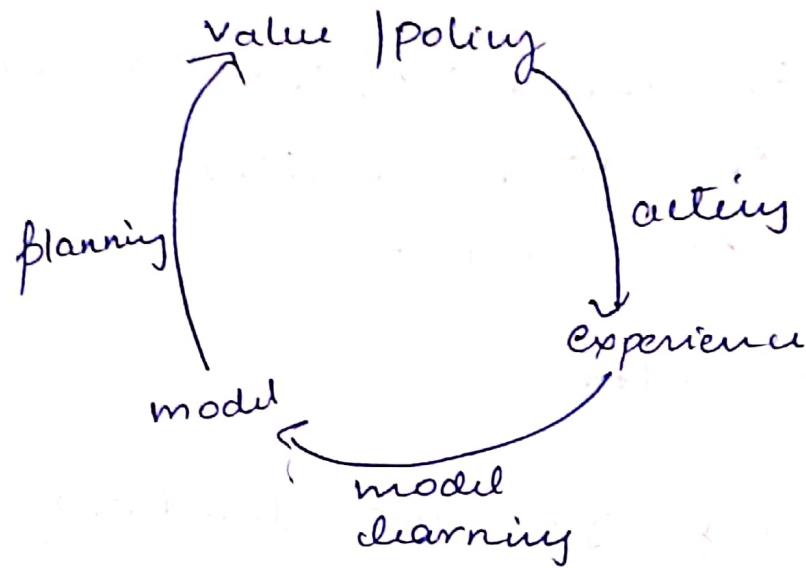
Natural policy gradient

- parameter independent

Integrating Learning

and planning.

- learn model from experience and use planning to construct value fn or policy.
- model tells about rewards and transition prob.



### Advantages:

- if learn value fn. or policy directly is hard (e.g. playing chess), model can provide more compact info.
- can learn using supervised learning
- reason about model uncertainty (choose to go areas where model is uncertain)

### disadvantages:

- first learn model, then value fn. (two sources of error)

### Model

- representation of a MDP ( $S, A, P, R$ ) parameterized by  $n$
- assume some  $S$  and  $A$  are known
- $M = (P_n, R_n) \rightarrow$  state transitions

$$P_n \approx P$$

$$R_m \approx R$$

$$S_{t+1} \sim P_n(S_{t+1}|S_t, A_t)$$

$$R_{t+1} \sim R_m(R_{t+1}|S_t, A_t)$$

Typically assume  $S$  and  $R$  are indep.

$$P[S_{t+1}, R_{t+1} | S_t, A_t] = P[S_{t+1} | S_t, A_t] P[R_{t+1} | S_t, A_t]$$

goal: estimate  $M_\pi$  from exp.  $\{S_1, A_1, R_2, \dots, S_T\}$

This is supervised

→ took action  $A_1$  from  $S_1$

$S_1, A_1 \rightarrow R_2, S_2 \rightarrow$  new state

↳ immediate reward

$S_2, A_2 \rightarrow R_3, S_3 \rightarrow$  previous

⋮

$S_{T-1}, A_{T-1} \rightarrow R_T, S_T$ .

$S, a \rightarrow r$  is a regression problem

$S, a \rightarrow s'$  is a density estimation problem.

Pick loss functions such as MSE, KLD etc.

## types of models

- Table lookup
- Linear expectation
- linear gaussian
- gaussian process
- Deep belief Network
- etc.

## table look up model

- count visits to each action - state pairs  
→ using this, build a probability distribution

or

- reward experiments seen so far and randomly pick up tuples of state - action

## Planning with a model

Solve MDP to find optimal value, policy to pick actions

### algor:

- Value iteration
- Policy iteration
- Tree search

etc

## Sample based planning.

- use the model only to generate samples
- sample exp from model
  - Then apply model-free RL to samples
  - MC control
  - sarsa
  - Q-learning
- often more efficient
- solve for the imagined experiences.

## Planning with an inaccurate model:

for an important model,

$$P_m, R_m \neq P, R$$

Model based RL is only good as the model that is learnt.

### Sols:

- when model is wrong, use model-free
- reason explicitly about model uncertainty  
(bayesian approach)

## Integrated architectures.

- two sources of experience

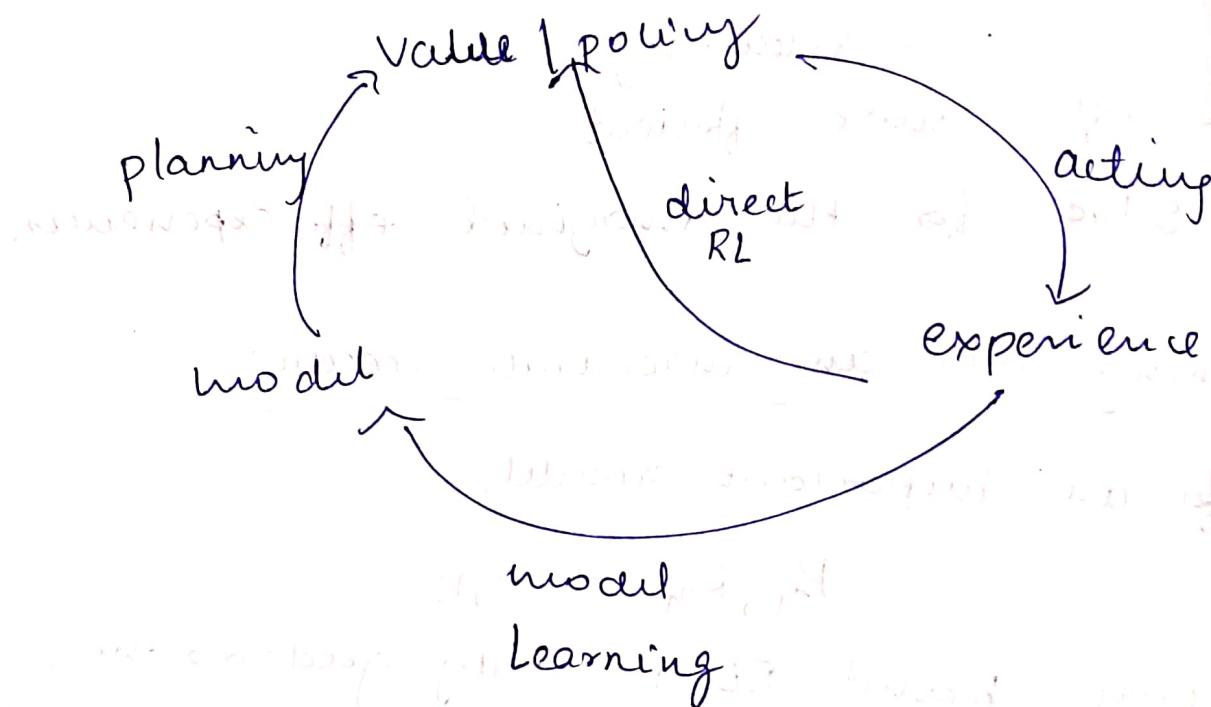
- real exp, sampled from envi  
 $s' \sim P_{s,s'}^a, R = R_s^a$  (true MDP)

- simulated exp, sampled from model  
(approx MDP)

$$\rightarrow (s' | eA) \quad R \sim R_m(R | s, A)$$

## Dyna

- learn a model from real experience.
- learn and plan value fn from real and simulated exp.



## Dyna - Q

Initiation  $Q(s, a)$  and  $\text{model}(s, a)$

- i)  $S \leftarrow$  current state
- ii)  $A \leftarrow \epsilon\text{-greedy}(s, Q)$
- iii) Execute  $A$ , observe  $R$ , state  $s'$  (see where we end up)
- iv)  $Q(s, A) \leftarrow Q(s, A) + \alpha [R + \gamma \max_a Q(s', a) - Q(s, A)]$   
[Normal Q learning]
- v)  $\text{Model}(s, a) \leftarrow R, s'$  update model using supervised L band on  $R$  and  $s'$
- vi) Repeat n times:  
-  $S \leftarrow$  random seen state and a action

- $R, S' \leftarrow \text{model}(S, A)$  Imagine the transition
- Q learning update, update Q value from when started to get Q value we ended up ( $s'$ )

## Simulation Based Search.

### forward search:

- these algorithms don't explore every state space
- They select the best action by lookahead.
- They build a search tree with current state  $S_t$  at the root
- Using a ~~MDP~~ model of MDP to lookahead
- No need to solve whole MDP, Only starting from now.

## Simulation-based Search.

- Forward search using samples
  - Apply Model for RL to samples
- $$\{S_t^k, A_t^k, R_{t+1}^k, \dots, S_T^k\}_{k=1}^K \sim H_n$$
- & samples.

Then,

MCC  $\rightarrow$  MC search

$\rightarrow$   $\rightarrow$  much

## Simple Monte Carlo Search

- Model  $M_v$  in simulation policy  $\pi$
- for each action  $a \in A$ 
  - Simulate  $k$  episodes from state  $s_t$
  - Evaluate actions by MC evaluation
- Select next action with max value

## Monte Carlo Tree Search (SOTA 2016) and in 60

- Model  $M_v$  in simulation policy  $\pi$
- Simulate  $k$  episodes from  $s_t$  using a current simulation policy  $\pi$
- Build a search tree based on all state-action pairs using MC evaluation
- After search is complete, pick actions with highest value in the root.
- After every simulation,  $\pi$  improves
- Each simulation can pick
  - Tree policy: pick actions to maximize  $Q(s, a)$
  - Default policy: pick actions randomly for  $s, a$  not in tree.
- Repeat this
  - Improve tree policy with G-greedy etc.
- Converges on optimal search tree.

## GO position evaluation

$R_T = 0$  for all non-terminal states

$$R_T = \begin{cases} 0 & \text{lose} \\ 1 & \text{win} \end{cases}$$

policy  $\pi = \{\pi_B, \pi_W\}$  for both players.

$$V_\pi(s) = E_\pi[R_T | S=s] = P[\text{win} | S=s]$$

$$V_t(s) = \max_{\pi_B} \min_{\pi_W} V_\pi(s) \rightarrow \text{tree band.}$$

## Advantage of MCTS

- highly selective best-first search
- evaluates states dynamically.
- uses samples to break C of dim.
- parallelizable

## TD-Search

- TD-search applies Sarsa to Sub-MDPs
- For each step in sim, update A, V by Sarsa
- Select actions ε-greedy
- May also use for approx for Q.

## Dyna-2

- maintains 2 value functions
  - like a long term memory from real experience.
  - and a short term memory as a search tree.

## Exploration and Exploitation

Exploitation: Make decision based on current info.

Exploration: gather more information by doing something else.

### Types of exploration

- Random exploration
  - explore random action e.g.:  $\epsilon$ -greedy.
- Optimism in the face of uncertainty
  - estimate uncertain in value
  - prefer to explore states/actions with highest uncertainty.
- Information State Space
  - consider agent's info as part of its state
  - lookahead to see how info helps reward.  
(computationally hard)

~~How~~ Our car explore in

- state space  $\rightarrow$  action space
- parameter space  
(parametrized policy  $\pi$ )  
(adv: consistent exploration)  
(disadv: doesn't know alt. states)

## Multi-Armed Bandit

- simplification of MDP.
- MAB is a tuple of  $\{A, R\}$
- $A$  is a set of Actions or arms
- $R^a(r) = P[R=r | A=a]$  is an unknown prob distribution over rewards.
- At every  $t$ , select one action, env generates a reward.
- maximise cumulative reward over time
- regret: opportunity loss for one step
- action-value: mean reward for action  $a$
- optimal-value: best action to choose

$$\textcircled{2} \quad q(a) = E[R|A=a]$$

$$\textcircled{3} \quad v_* = q(a^*) = \max_{a \in A} q(a)$$

$$\textcircled{1} \quad l_t = E[v_* - q(A_t)]$$

↓ optimal      ↗ payoff of picked action

- total regret

$$L_T = E \left[ \sum_{t=1}^T v_* - q(A_t) \right]$$

- Maximum cumulative reward  $\Rightarrow$

minimum total regret

- Count  $N_t(a) \rightarrow$  expected no. of selection for  $a$

- gap  $\Delta_a$  is  $v_* - q(a)$

$$\therefore L_T = E \left[ \sum_{t=1}^T v_* - q(A_t) \right]$$

$$= \left[ \sum_{a \in A} E[N_t(a)] \right] (v_* - q(a))$$

$$\sum_{a \in A} E[N_t(a)] \Delta_a$$

Regret is a function of counts and gaps.

- Every time gap  $\Delta_a$  is large, make sum count is low. If gap is low, let count be more.

- problem: gap are not known, because we don't know  $v_*$

greedy algo with this

- greedy algo can pick sub-optimal action forever.
- it has linear total regret.

Algorithms with Optimistic Initialisation

- really good idea
- initialising values to max possible  $Q(a)$ ,  $v_{max}$
- Then act greedily

$$A_t = \arg \max_{a \in A} Q_t(a)$$

- Encourages exploration of things not known.
- Linear total regret.
- Some unlucky samples can lock out some states forever.

$\epsilon$ -greedy

- explores forever
- Constant  $\epsilon$  ensures min regret
- linear total regret.

decaying  $\epsilon$ -greedy

- decay  $\epsilon$ .
- logarithmic regret.
- this needs advance knowledge of gaps.

- Total regret is at least  $\log n$  in steps.

### optimism in face of uncertainty

- pick more uncertain actions

### UCB1

$$A_t = \underset{a \in A}{\operatorname{argmax}} \left[ Q_t(a) + \sqrt{\frac{2 \log t}{N_t(a)}} \right]$$

total time steps  
no. of times of this action

M C estimate of Q value

This achieves  $\log n$  total regret

### Bayesian Bandits

- exploit prior distn on  $[R]$
- Use posterior to guide exploration

### Probability matching

- Selects action according to prob that it is the best one.
- Can be difficult to compute  $\pi(a)$  analytically from posterior.

### Thompson Sampling

- Sample based probability matching

$$\pi_t(a) = E [I(Q_t(a) = \max Q_t(a')) | R_1, \dots, R_{t-1}]$$

- Every step, pick sample from posterior
- Then pick the one which is best

Then algo keeps going forever (can't impose budget)

### Info-state space

- Info gain is high in uncertain situations
- so if we know value of info, we can trade off exploration & exploitation optimally.
- see bandit problem as MDP instead of one-step
- At each state, there is an info state  $\tilde{S}$  summarising all info seen.
- An action  $A$  causes a transition to a new info state  $\tilde{S}'$ , with prob  $\tilde{P}_{\tilde{S}\tilde{S}'}^A$
- So there is now an MDP  $\tilde{M}$  in the augmented info state space

$$\tilde{M} = \{\tilde{S}, A, \tilde{P}, R, \gamma\}$$

- Consider Bernoulli Bandit (like coin flip).

- $R^a = B(\mu_a)$   $\mu_a \rightarrow$  prob.
- find highest  $\mu_a$ .
- Info state  $\tilde{S} = \{\alpha, \beta\}$

It can be

- $\alpha$  counts where reward = 0
- $\beta$  counts where reward = 1
- To solve this,

- start with a  $\text{Beta}(\alpha_a, \beta_a)$  prior over  $\alpha_a$
- update posterior for new action
  - $\text{Beta}(\alpha_a + 1, \beta_a)$  for  $r^0$
  - $\text{Beta}(\alpha_a, \beta_a + 1)$  for  $r^1$
- this is Bayes adaptive MDP

## Contextual Bandits

- add states.
- $S = P(S)$ , prob distn over states.

## Other approaches

- reconstruct model
- for unknown states, give  $r_{\max}$
- Be very optimistic towards uncertainty
- solve MDP using any algo.

## game-playing

Best response is other player(s) become part of environment. This is just an MDP now.

NE is when every player has found best response to everyone else.

## Perfect info games

- all players can access info of others

## Imperfect info games

e.g. poker

### Value fn:

$$\text{policy } \Pi = \{\pi^1, \pi^2\}$$

$$V(s) = E_{\Pi} [u_t | s_t = s]$$

A minimax fn maximises own's expected return  
& app minimises opponent's expected return