

## 70 Policy-based methods.

↓  
RL is ultimately about learning an optimal policy from interaction with the environment

↓  
till now, saw value based methods, where we need to estimate an optimal action value function, to estimate an optimal policy.

but for large state space environments

← okay for env. with small state spaces.

Some sort of approximation has been formulated till now.

→ allows non-linear function approximation  
learn to represent the action value function in the form of a neural network

↓  
discretization

↓  
function approximation

↓  
tile coding

radial basis function

↓  
env. state as i/p to neural network

↓  
o/p was value of each possible action

↓  
and as per the previous method,

action is selected on the basis of these values, which can done either on the basis of greedy or  $\epsilon$ -greedy.

(2)

① table

↓  
neural net (action values)

↓  
choose based on preferred technique.

next process can stay same.

\* But optimal action value function has been estimated before estimation of optimal policy.



can policy be found directly? without worrying about the value function at all.



yes → policy based methods.

→ neural network to approximate policy → state as i/p.



greedy ↓ deterministic policy  
← in the o/p layer, rather than return action values, it can return probabilities of the actions, & accordingly sample from those. → sampling: stochastic policy.



i.e. 2 possible actions, 2 nodes in the final layer.



agent uses this policy to interact with the env. by passing the state to the network.



we need to then select appropriate weights for the network, so that for each state that is passed, o/p is action probabilities where the optimal action is most likely to be selected.



initially random values and then fine-tuned to maximize reward as more interaction takes place, gradually mastering the task.

72 \* for continuous action space.

↓  
neural network has one node for each  
action entry (or index)  
(i.e. for each constituent of  
the action vector)

↓  
for example: bipedal walker

NUM	NAME	min	max
0	Hip-1 ( $z/v$ )	-1	1
1	Knee-1 ( $z/v$ )	-1	1
2	Hip-2 ( $z/v$ )	-1	1
3	Knee-2 ( $z/v$ )	-1	1

↓  
action: vector of four numbers, so the  
o/p layer of the policy network will have  
four nodes

↓  
since every value must be a number between  
-1 and 1, a tanh activation will be added.

↓  
similarly for continuous mountain car  
benchmark

\* Hill climbing: optimization technique used to find  
maximum of a function

agent's goal: to maximize the expected return  
for any environment,  $\pi$   
weights in the network:  $\theta$



mathematical relationship b/w  $\theta$  & expected return  $J$

(73)

as  $\theta \rightarrow$  encodes the policy  $\rightarrow$  determines the action  
(affects the policy)

$\downarrow$   
influences reward

$\downarrow$   
used to get expected return,  $J$

$$J(\theta) = \sum_z P(z; \theta) R(z) \rightarrow \text{it exists}$$

(mathematical relation b/w  $\theta$ , and  $J$ )

$\downarrow$   
goal is to find  $\theta$ , that maximize  $J$

$\downarrow$   
considers a case where the neural net has only 2 weights:  $\theta_1, \theta_2$

$\downarrow$   
estimating  $\theta_1, \theta_2$  values to maximize the reward, can be done using technique called gradient ascent.

i.e. rather than moving in the direction opp. to the gradient (as done in gradient descent) to find minima, we move in direction of the gradient.

$\downarrow$   
begins with an initial guess.

$\downarrow$   
this process is repeated

$\downarrow$   
to do this, gradients need to be evaluated. (later)

$\downarrow$   
Now: hill climbing

74) As with gradient ascent, initiate with a random set of values for  $\theta$ . collecting an episode corresponding to the policy for the weights of the network, and then the record is returned.

↓  
add random noise to weights <sup>gaussian or some other model</sup> another set of candidate values  
↓  
if better good. we change the current weights to the new ones,  
↓  
else we go back to our current net weights.

↓  
repeated till optimal policy not approximated.

hill climbing may lead to being stuck in a local minima.

(\*) pseudocode : Hill climbing :

- i) initialize the weights  $\theta$  in the policy arbitrarily
- ii) collect an episode with  $\theta$ , and record the return  $g$ .

$$\theta_{best} \leftarrow \theta$$

$$g_{best} \leftarrow g$$

iii) Repeat until env. solved :

\* Random noise ~~added~~ added to  $\theta_{best}$ , to get a new set of weights.  $\theta_{new}$ .

\* collect an episode with  $\theta_{new}$ , to get  $g_{new}$

\* If  $g_{new} > g_{best}$  :

$$\theta_{best} \leftarrow \theta_{new}, \quad g_{best} \leftarrow g_{new}.$$

(75)  
\* general class of approaches, that find  $\arg \max_{\theta} J(\theta)$   
through randomly perturbing the most recent best estimate  
as stochastic policy search. Likewise  $J$ , can be referred to as  
an objective function. (we would like to maximize it)

↓  
one improvement could be to select a number  
of neighbouring policies at each iteration and then  
picking the best among them.

↓  
gives the agent an idea about the neighbourhood of  
the policy.

↓  
steepest ascent hill climbing

(reduces the risk of selecting a  
next policy that may lead to a suboptimal  
solution)

• getting stuck in local minimas → a) random restarts

b) simulated annealing

↓  
predefined schedule to  
control how the policy  
space is explored

↓  
initially : large noise parameter (broad neighbourhood  
to explore)

↓  
gradually noise is reduced, as  
we get closer & closer to the optimal selection

↓  
analogous to annealing iron : heating it up  
and letting it cool down  
gradually.



76) \* adaptive noise

↓

as better policy is estimated, radius for search can be reduced. for generating the next policy

↓

translates to reducing or decaying the variance of the Gaussian noise we add,

but if don't find a better policy, it is better to increase the search radius, and to continue exploring from the current best policy.

↓

reduces the probability of it getting stuck.

→ Black-box optimization:

↓

techniques used so far, do not have any idea about their use case, and are as such not dependent on the use case. These algorithms just care about approximating the  $\theta$  that maximize  $J$ , and can be applied to any similar case.

↓

other black box optimization techniques:

i) cross-entropy method

ii) evolution strategies.

i) cross-entropy method:

hill climbing  
→ steepest ascent hill climbing

leveraging useful information from the weights that aren't selected as the best

↓

rather than selecting the best, selecting the top 10 or 20%, and then taking its average.

ii) evolution strategies. <sup>→ name originally comes from biological evolution</sup>

↓  
best policy : weighted sum of all candidate policies,  
where policies with higher return has higher say or get  
a higher weight.

↓  
most successful individuals have more  
influence on the next generation.

★ Why policy-based methods?

- i) simplicity
- ii) stochastic policies.
- iii) continuous action space.

→ directly getting to the problem  
at hand, rather than  
worrying on computation of  
some intermediate step.

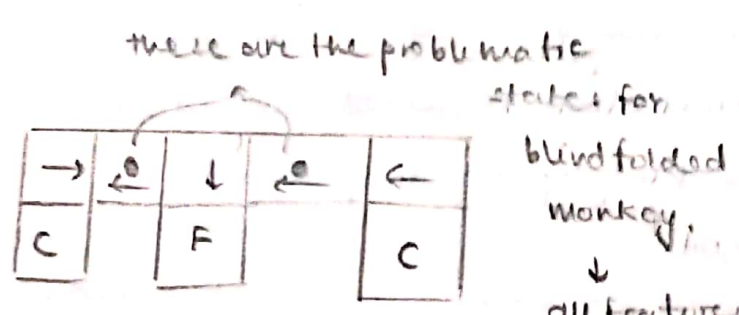
(well suited ;  
(value based isn't)

\* being able to get a  
true stochastic policy.

↓  
also avoids storing additional data,  
that may or may not be useful

↓  
ε-greedy is sort of a hack

aliased states :: 2 or more states that are perceived to be  
identical, but are actually different.



↓  
partial observability  
of the environment

↓  
all features are identical

↓  
action values are also equal

← Since they map to the same state,  
∴ corresponding best action will  
become

could be wrong for  
at least 1  
aliased state

↓  
oscillations

↳ he could come out using ε-greedy



(78) but this can be really inefficient & could take a lot of time.

↓

and high epsilon could lead to bad actions.

↓

equal prob to both actions!

↓

value based approach tends to learn a deterministic or near-deterministic policy, whereas a policy based approach can learn the desired stochastic policy.

Finding the most profitable or rewarding action over a continuous action space is an optimization problem in itself.

↓

(the more features the action vector has, harder it gets to find the other corresponding maxima)

↓

hence it is better that we can map a given state to an action directly.

↓

even ~~then~~ if the resulting policy is a bit more complex, computation required would drop significantly, enabled by a policy based method.