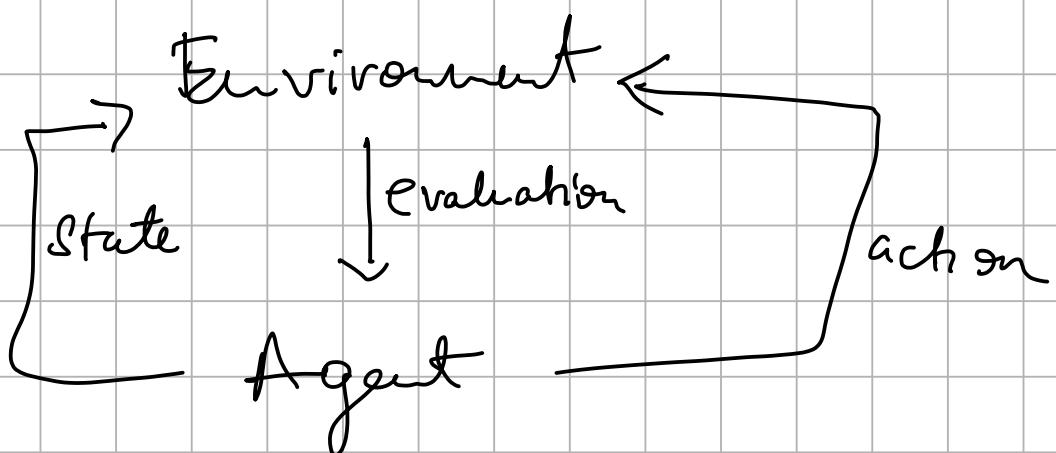


# Reinforcement Learning

RL framework



- Learn from close interaction
- Stochastic environment
- Noisy delayed scalar evaluation
- maximize a measure of long term performance

# Temporal Difference

- Simple rule to explain complex behaviors

- Intuition :

Prediction of outcome at time  $t + 1$  is better than the prediction at time  $t$ . Hence use the later prediction to adjust the earlier predictions.

- Has had profound impact in behavioral psychology and neuroscience

# Explore Exploit Dilemma

- One key question - the dilemma between exploration and exploitation.
- Explore to find profitable actions
- Exploit to act according to the best observations already made
- Bandit problems encapsulates 'Explore vs Exploit'

# Immediate Reinforcement Learning Problems

- Problems when:  
at every time  $t$ , we pick  
an action  $a_t$  and get the  
reward  $r_t$
- We need exploration for any  
reinforcement learning problems
- We must stop our exploration and  
'exploit' a solution that was  
found (we cannot keep exploring  
forever)

Set of actions  $A = \{1, 2, \dots, n\}$

and each action gets a pay-off or cost or reward or evaluation

note

One action can have multiple rewards - Example:

Tossing a coin  $\leftarrow$  action

Reward  $\leftarrow$   $\begin{cases} 1 & \text{if head} \\ 0 & \text{if tails} \end{cases}$

So different reward for the same action

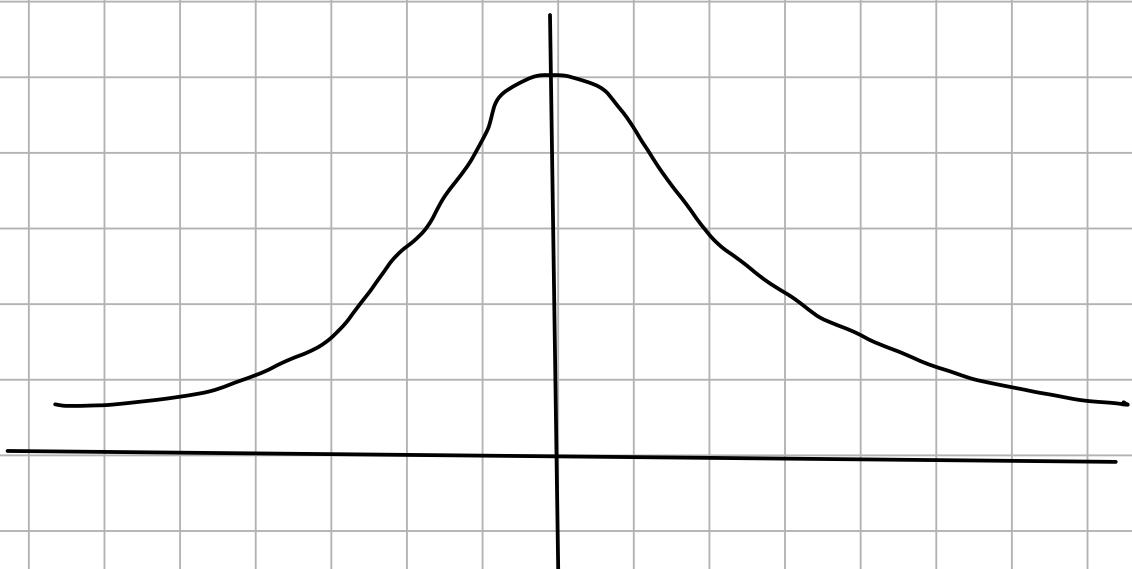
So the reward depends on the probability

So the coin toss reward is

Sampled from a Bernoulli Distribution

Let's say any action form

$A = \{1, 2, \dots, n\}$  Samples from Gaussian distribution



$q_{\pi}(a)$   
↑ action

gives the best reward of action  $a$   
(true expected reward)

# Bandit Problems [multi arm bandit]

note

actions are called arms sometimes

Bandit Problem means that there are  $n$  number of actions and we get a reward after each action taken

Solving multi arm bandit problems.

one way is

- Asymptotic Correctness

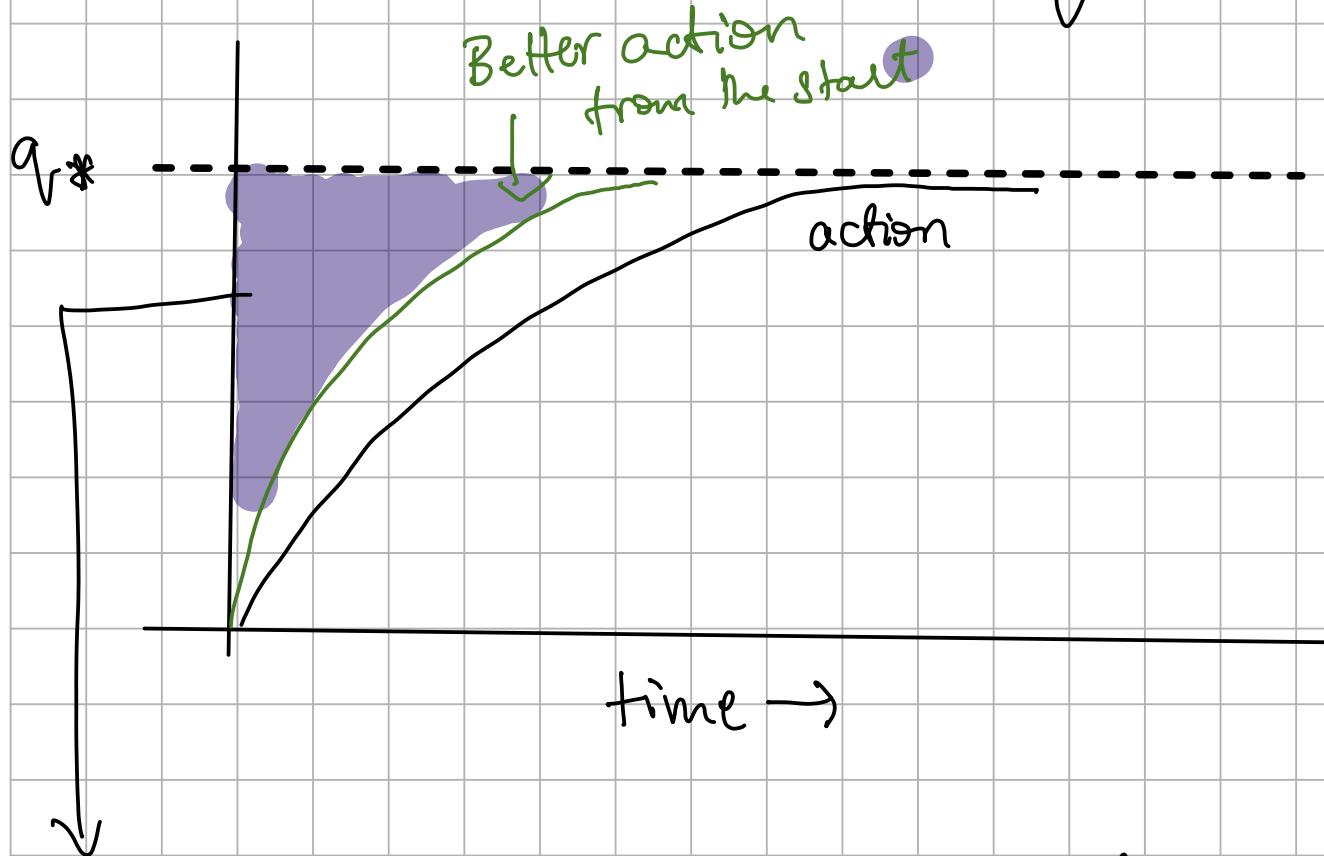
L this is defined as we keep taking

actions , but eventually we must select

the 'arm' that has the highest reward

(how quickly can we get that arm is another issue separately)

- Regret Optimality
  - ↳ trying to maximize the total reward from the beginning



This area is or loss, this is called regret  
 So in regret optimality, we are  
 trying to decrease the regret as much  
 as possible

- PAC Optimality

Probably correct  $\leftarrow$  can be correct or wrong

Approximately correct  $\leftarrow$  it is definitely wrong but close to the correct answer

PAC  $\leftarrow$  Probably approximately correct (highly)

PAC is represented as  $(\epsilon, \delta)$  pac

where  $\epsilon \leftarrow$  approximability variable,

$\delta \leftarrow$  probability path variable

So for formulation

$$\text{Prob} [q_*(a) \geq q_*(a^*) - \epsilon] \geq (1 - \delta)$$

$a^* \leftarrow$  action being taken