

31st July '25

Reinforcement Learning

valuation for a particular state for future to get the greatest amount - Value function.

pure method - trial and error method (most of them).

expecting value of each and every state of (probability) methods.

model based methods are also used

→ evolutionary methods - hill climbing method

(searching in space of probability)

can also be used to solve reinforcement learning

if you can solve the problem it would be more optimal solution for RL.

→ regularly used

in evolutionary methods, small space is given

- large branching factor (10^{22})

- stochastic environment

In the above case, it is more effective

e.g. Tic Tac Toe game - 9 initial states, moving forward reduces.

↳ learn agent through reinforcement learning

classical methods of game theory

↳ minimax algorithm - player is imperfect (considering)

↳ search space discloses the narrow area to learn the model but it cannot learn efficiently

agent could be in draw).

cannot give you up-to-date solution.

not optimal

classical approaches : dynamic approach (required prior information)

for any RL problem (no prior information is given)

it is needed to gather information from experience

→ branching factor in this game is comparatively less so evolutionary methods can be used.

losing the game → reward 0
winning the game → reward 1.
probability for all states → 0.1
subsequently change the policy (reward is changing at the time of interaction with the environment).

→ one state to another state movement with greatest reward but sometimes we also need to explore (may be not having the greatest reward) (exploratory action).

dilemma of exploration & exploitation.

(killing the probability, closing the regime)
 $\begin{array}{c} A \\ \text{---} \\ B \end{array}$ 0.6 0.3

so, move to the state that are close to that probability

→ A's prob > B's prob. → decrease the probability of state A & viceversa.

(Updating parameter, terms of learning) → for refining the environment

(unsupervised learning - choosing the gradient)

but in RL → yourself find out the next state given the environment

the above is the.

Simpler approach of the reinforcement learning

(function approximator) developing through a supervised learning

(using a function approximator to find out the percentage of success and failure of the policy or

4th Aug/25

unsupervised learning

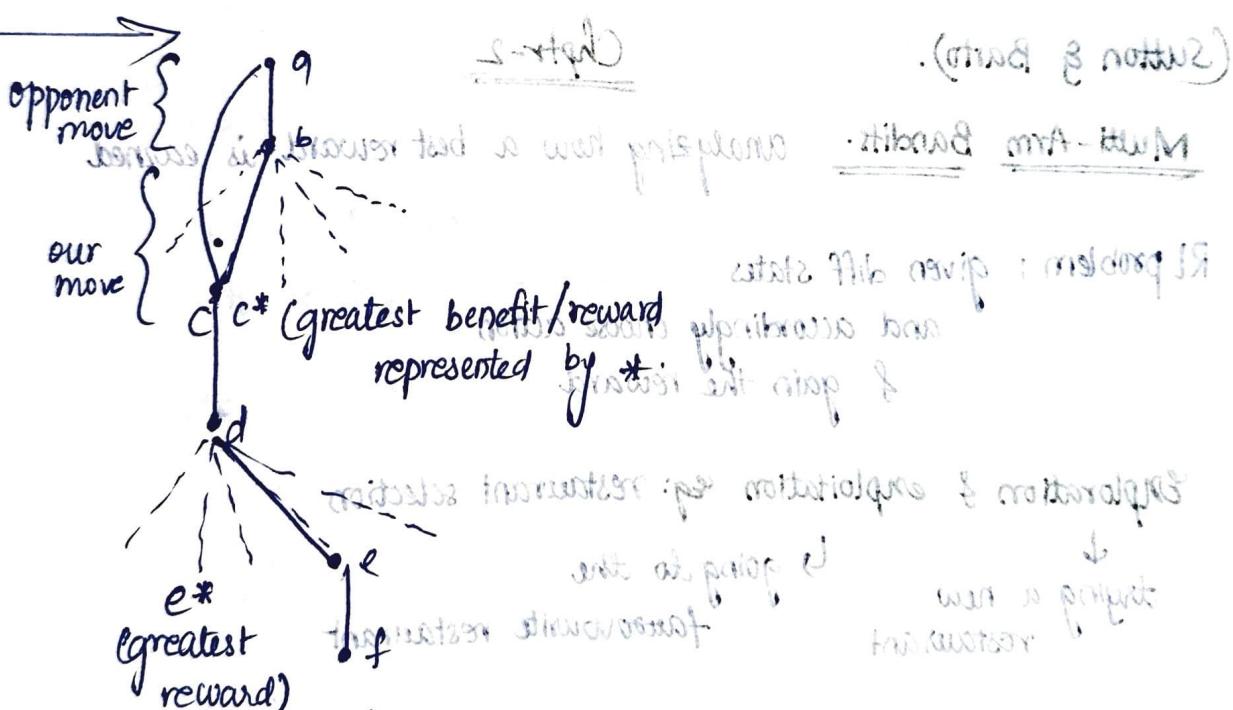
not unsupervised learning
↳ finding pattern

→ here it is not the primary goal

evolutionary methods → can't be optimal for large searching space

Tic-tac-toe : halite & valves

minimax } classical techniques
DP → creating model of opponents moves.
evolutionary methods.
RL method



→ temporal difference learning method

$$V(s) \leftarrow V(s) + \alpha [V(s') - V(s)]$$

state before greedy move

small +ve state after

(step-size parameter)

↳ influences rate of learning

Temporal difference:

$t+1$ prediction is better than prediction at t ? Intuition: missing problem?

Explore - Exploit dilemma:

(no supervision before so, we encapsulate this)

Simplest version of RL \rightarrow Bandit problems

problems.

encapsulate

explore & exploit

Slot machine \rightarrow collecting money.

\downarrow few may hit the jackpot (or) nothing.

can be a type of Bandit

(Sutton & Barto). — book name

Chptr-2

Multi-Arm Bandits. analyzing how a best reward is earned

RL problem: given diff states

and accordingly choose action

& gain the reward

Exploration & exploitation e.g: restaurant selection

\downarrow
trying a new
restaurant

\hookrightarrow going to the
favourite restaurant

$q_* \rightarrow$ true expected value after performing an action

Action: 88% chance of getting 0

12% chance of getting 100

} for this action,

$$q_*^{(2)} = 0.88 \times 0 + 0.12 \times 100 = 12.$$

(true value)

randomly - 10% & 85% equiprobable. — Action (5).

$$\text{Number line from } -10 \text{ to } 35 \text{ with tick marks at } 0, 10, 20, 30. \text{ The interval } [-10, 35] \text{ is shaded in blue.}$$

deterministic \rightarrow best action.

deterministic \rightarrow best action.
stochasticity \rightarrow reduced (the chances are reduced to lose)
gets more real. ultimately to

k-armed Bandit problem.

armed Bandit problem. ~~Goal~~ - ~~by~~ rewards are to
distribution is unknown - so, knowing them & rewards are to
not be exploited - explore first

Exploration / Exploitation dilemma:

Exploration dilemma: 

$A_t = A^* t \rightarrow$ exploiting } mathematically.

$A^t \neq A^{*t}$ (not) exploring

→ Asymptotic correction: no bounded condition.

for algorithms in every moment making an 'optimal' action to solve RL.

any is asymptotic correctness of the problem

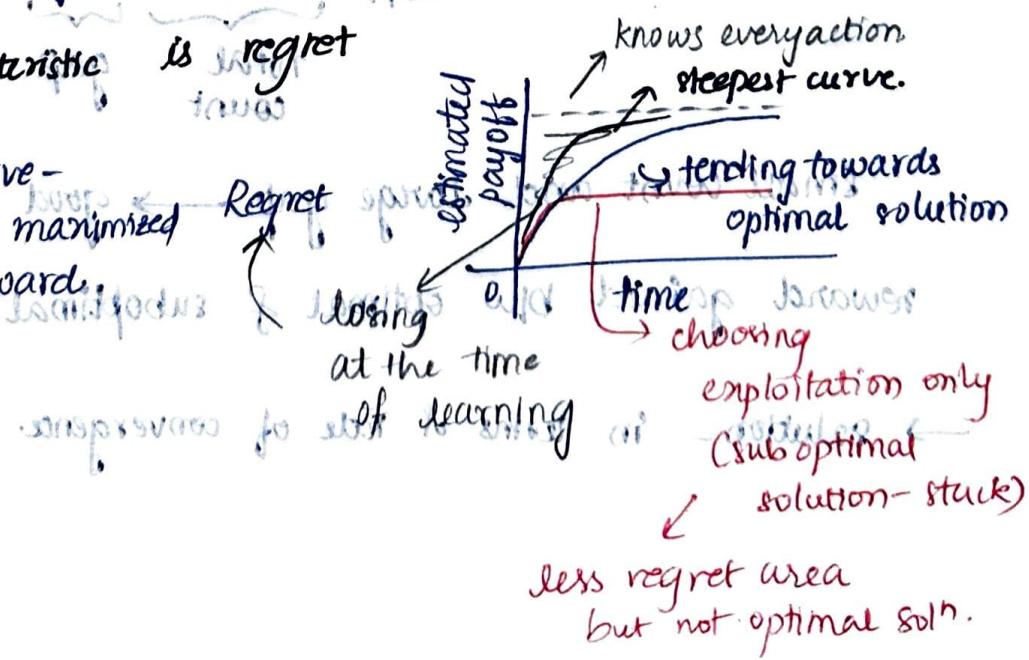
(guaranteed) in future)

→ Another characteristic is regret

the steepest curve -

gives the maximized reward

loop
wonder,
handed is not
mention



Regret:

(a) most A = suboptimal actions & b = good actions

6th Aug's:

$$V^* = \max_{a \in A} q^*(a)$$

estimated
optimal value.

(and other actions will be better → - regret)

Regret → opportunity loss for one step

$$l_t = E[V^* - Q(a_t)]$$

total regret = total opportunity loss.

$$l_t = E \left[\sum_{r=1}^t V^* - Q(a_r) \right]$$

↳ timestamp = 1

→ how many times we are considering an action - count

$$\Delta a = \underbrace{V^* - Q(a)}_{\text{estimated value of action } a}$$

$$l_t = E \left[\sum_{r=1}^t V^* - Q(a_r) \right]$$

$$= \sum_{a \in A} E[N_t(a)](V^* - Q(a))$$

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

total Δ gap / total count

small count and large gap → good algorithm.

reward gained b/w optimal & suboptimal action.

→ solution - in terms of rate of convergence

↓
good
solution
for a bandit
problem

→ forever exploring & never exploring gives Linear regret
 ϵ -greedy (epsilon greedy algorithm) - balances exploiting &
decaying ϵ -greed → explores with ϵ & explores

Solution:

- 1) Asymptotic correctness
- 2) regret optimality.

3) ~~pac optimality~~ probably approximately correct

↳ ϵ -pac optimality (Epsilon delta)

(confidence)

hypothesis: close to optimal
with gap ϵ

probability of correctness

* $Q^*(a) - \text{expected/best true value} > Q^*(a^*) - \epsilon$

↑ speak of between target dist. and next
(ϵ -gap)

probability of this occurrence should be greater than $(1-\delta)$

if this is guaranteed solution is having the pac optimality

still you need to behaviors in \mathcal{A}) bottom phas... nalgap ←

⇒ no solution - log theory.
(.loglog)

Complexity of regret helps one to minimize regret.

& partial regret \rightarrow total regret should always be greater than log t

where t is the time stamp.

$$\lim_{t \rightarrow \infty} \Delta_t \geq \log t \sum_{\Delta a > 0} \frac{\Delta a}{KL(R_a || R_a^*)}$$

(partial regret)

Hard problems have similar looking with different means.

(Action regret) (difference in q distributions)

according to plaidor

$\Delta a \rightarrow$ gap b/w optimal & other arms + size of gap

$$\sum_{\Delta a > 0} \frac{\Delta a}{KL(R_a || R_a^*)} \rightarrow \text{if this is small, then the total regret would be large.}$$

(cont) best action & reward outcome bit is plaidor

Action-value methods:

→ Epsilon-greedy method (ϵ is considered in such way those 90% exploit 10% explore many times the actions are explored)

disadvantage - so much computation & memory usage
coz keeping track of all records & values would take large time & storage.

→ Incremental implementation.

→ stationary vs non stationary environment

distribution of action is not changed	distribution of action is changed
(fixed)	

→ Optimistic initial values ~~initial state - it will be explored~~
1st requirement to solve bandit problem - exploration & exploitation
timestamp 0 → initialized every action = 0.
updating according to every action.
expected value > optimal value (explore more to reach optimality)

UCB action selection method:

considering uncertainty → not done enough sampling for that
would not contribute to an particular action being popular

Gradient bandit algorithm:

finding preferences of the action. → function

Chptr2 : Multi-Arm Bandit

www.india

27th Aug' 05

Associative search: Considering on the basis of recommendations & taking the actions that are required.

Episode: multiple time stamps.

(Ville de Québec et Sainte-Foy) avec un peu de la baie de Québec.

ACTION - VALUE METHODS:

- ↳ stationary environment, stable for \leftarrow plausible principles
 - ↳ average of past rewards in terms of estimation of the rewards.

$$Q_t(a) = \frac{\text{sum of rewards when } a \text{ taken prior to } t}{\text{number of times } a \text{ taken prior to } t}$$

$$\dim \mathcal{Q}_t(a) = q_*(a)$$

$N_t(a) \rightarrow \infty$ (converging to true values)

E-GREEDY ACTION SELECTION:

ϵ - value of epsilon is considered such a way to balance exploitation & exploration.

egj:

-A - one particular action.

In this, simple bandit problem

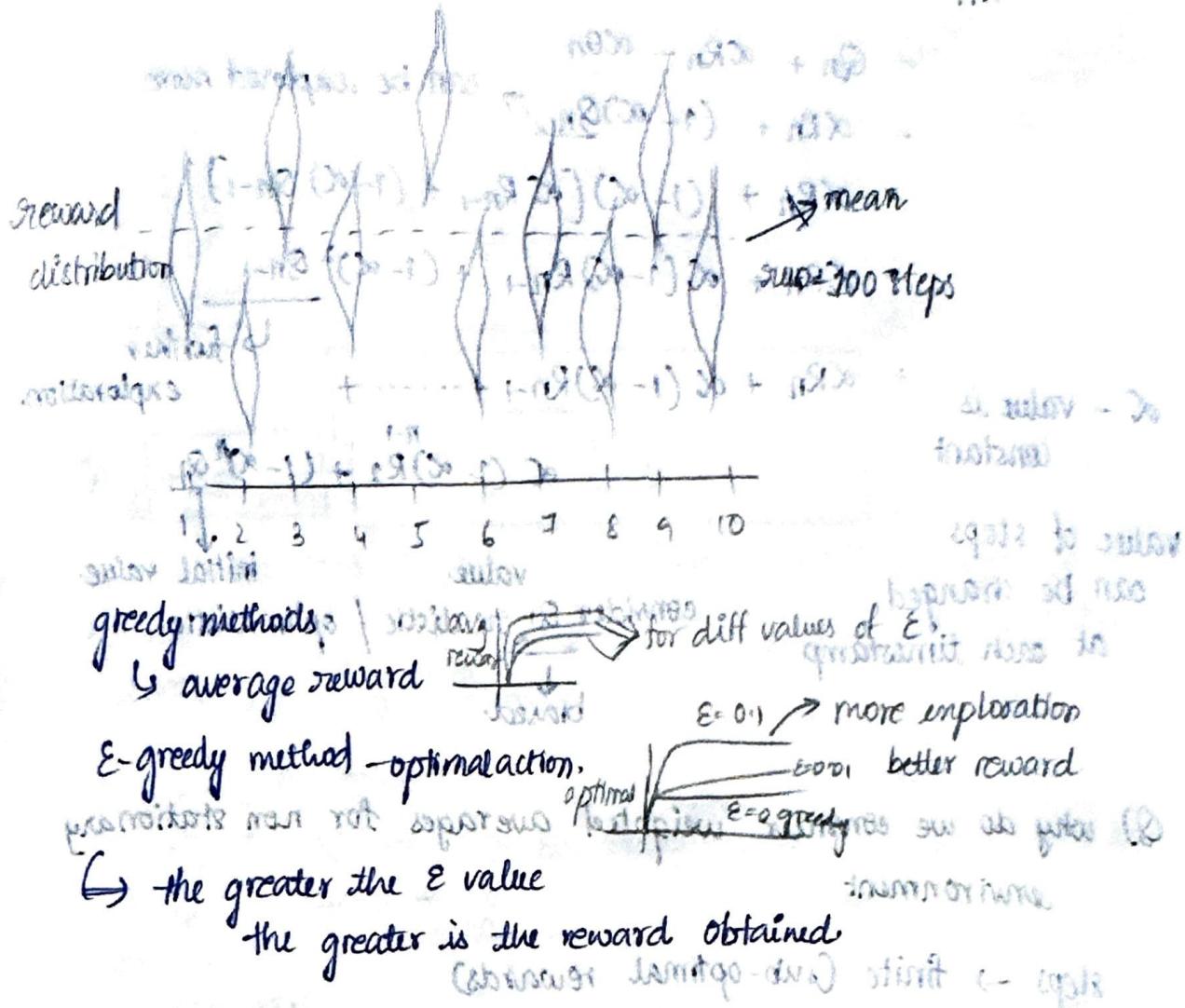
We consider only one particular action at a time.

$$Q(A) \leftarrow Q(A) + \frac{1}{N(A)} [R - Q(A)]$$

(increasing & checking other actions)
updating values using
incremental update formula
no. of times reward - old estimate value

10- Armed Testband

$$(m_1 - m_2) \cdot 20 + m_2 = 100$$



Tracking a Non-stationary problem:

changes made:

→ Over time true value of every action is changed

considering weighted averages.

$$Q_{n+1} = \alpha Q_n + (1-\alpha) R_i + \sum_{i=1}^n \alpha(1-\alpha)^{n-i} R_i$$

initial rewards $\cancel{\text{initial and rewards}}$ with all the rewards

initial action $\cancel{\text{initial and rewards}}$ with all the rewards

weighted averages $\cancel{\text{initial and rewards}}$ with all the rewards

→ more is the timestamp the more is its weight

i.e. Q_1 's weight < Q_{10} 's weight

$$Q_{n+1} = Q_n + \alpha [R_n - Q_n]$$

$$= Q_n + \alpha R_n - \alpha Q_n$$

$$= \alpha R_n + (1-\alpha) Q_n \rightarrow \text{can be explored more}$$

$$= \alpha R_n + (1-\alpha) [\alpha R_{n-1} + (1-\alpha) Q_{n-1}]$$

$$= \alpha R_n + \alpha (1-\alpha) R_{n-1} + (1-\alpha)^2 Q_{n-1}$$

$$= \boxed{\alpha R_n + \alpha (1-\alpha) R_{n-1} + \dots + \underbrace{\alpha (1-\alpha)^{n-1} R_1}_{\text{till } (n-1)} + (1-\alpha)^n Q_1}$$

↳ further exploration

α - value is constant

value of steps can be changed at each timestamp

consider Q_1 realistic / optimistic step initial value.
 ↓ biased.

Q) why do we consider weighted averages for non stationary environment

steps → finite (sub-optimal rewards)

$$\sum_{n=1}^{\infty} \alpha_n(a) = \infty \Rightarrow \text{converging to optimal reward}$$

$$\sum_{n=1}^{\infty} \alpha_n^2(a) < \infty \Rightarrow \text{estimation + noise is incorporated}$$

overtime noise is reduced.

shrink the npicy part

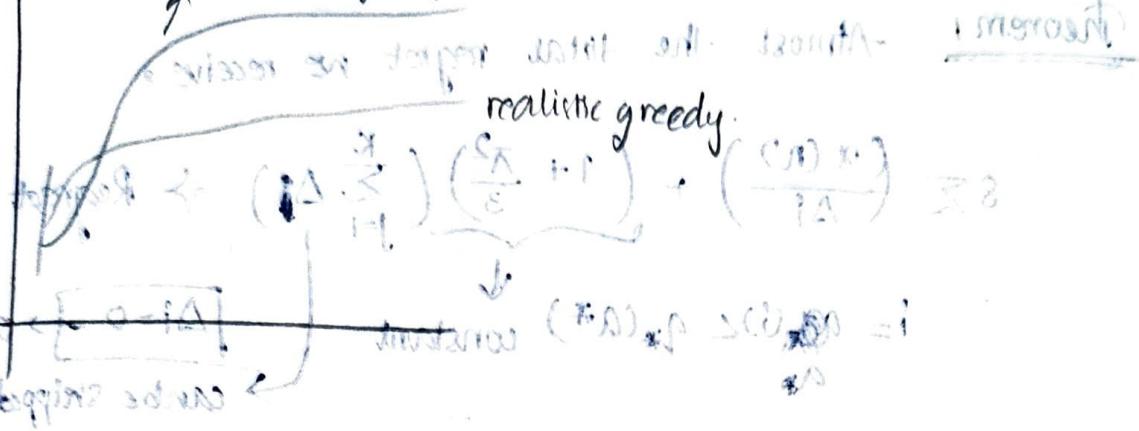
to do this condition should be satisfied

positive negative

$$\alpha_n = \frac{1}{n} \text{ and not } \frac{1}{n^2} \rightarrow \text{here it is ok only}$$

graph can be both negative and positive values so not considered to be normal if positive values.

optimistic greedy

18th August '25

$\text{Optimal solution, exploration of all arms is required of considering the uncertainty level of every arm}$

minimizing due to uncertainty level

$$\hat{Q}_j =$$

$$\sqrt{\frac{2 \ln n}{n_j}}$$

how many number of times we score equal unit in step n_j \rightarrow uncertainty level

estimation of each arm

$$[Q(1), Q(2)] \rightarrow$$

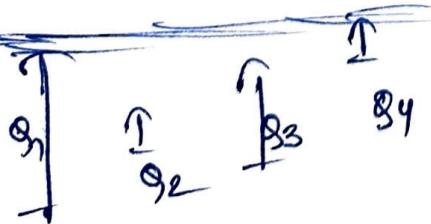
highest reward ??

(choose dominated action based on $Q(1) < Q(2) < Q(3) < Q(4)$ w/o uncertainty)

more and more moves left \rightarrow probability of getting very less compared to Q_4

totaling to expectation of each arm

as sampling increasing the confidence with more exploration of an action. $Q_1 = Q_2 = Q_3 = Q_4$



high

reward: $(Q_1)_{\text{act}} = \frac{1}{4} \cdot \frac{1}{4} \cdot \frac{1}{4} \cdot \frac{1}{4} = \frac{1}{256} \Rightarrow [Q(1)] = \frac{1}{256}$

increased the probability to obtain maximum

wants at still below Q_1

$$\Rightarrow [Q(1)] = \frac{1}{256}$$

UCB algorithm,

Theorem 1

Almost the total regret we receive.

$$8 \sum \left(\frac{a(n)}{\Delta i} \right) + \left(1 + \frac{\pi^2}{3} \right) \left(\sum_{j=1}^K \Delta j \right) \rightarrow \text{Regret}$$

$$i = q_*(i) - q_*(a^*) \downarrow \text{constant}$$

$\Delta i = 0 \rightarrow \text{optimal}$
can be skipped if \uparrow

i (arm)

sub optimal arm:

adjusting to decide learning

incurred some loss

while performing.

$$\Delta i = q_*(a^*) - q_*(i)$$

the sum of max values
is always at most 1.0
max prev. to small

Terms used in theorem:

$T_i(n) \rightarrow$ no. of times arm i has been up to n timestamps.

regret $n \rightarrow \sum E[T_i(n)] \Delta i$

(difference b/w expected value f. suboptimal action)

$x_{i,n} \rightarrow$ Random variable for reward that has been obtained for playing action i at n th timestamp.

$$E(x_{i,n}) = q_*(i)$$

expectation

of that random variable.

\rightarrow we would like to show.

$$E[T_i(n)] \leq \frac{8}{\Delta i^2} \ln(n) + \text{constant}$$

Concentration bound:

creating some bound around a variable which should not be exceeded.

→ Chernoff concentration bound:

the probability for expected value we achieved is not too far from true value.

↪ dependent of timestamps of less value we are considering.

→ Let x_1, \dots, x_n be a RV. with common $[0, 1]$ such that

$$E[x_t | x_1, \dots, x_{t-1}] = \mu \quad \begin{array}{l} \text{stationary problem} \\ \text{expected reward} \end{array}$$

Let $S_n = \underbrace{x_1 + \dots + x_n}_{n \text{ summands}} + \epsilon \geq 0$ stationary prove if $\epsilon \leq 0.02$

↪ To solve bandit problem with stationary bandit problem considers \rightarrow stationary environment

According to this bound, one particular action t doesn't cross a particular value.

$$\Pr[S_n \geq \mu + \epsilon] \leq e^{-\frac{\epsilon^2 n}{2\mu + 2\epsilon}} \quad \left. \begin{array}{l} \text{value of } \epsilon = 0.1 \text{ (there are 2 cases)} \\ \text{value of } \epsilon = 0.01 \end{array} \right.$$

As, with ϵ will increase
↑ $\epsilon \rightarrow P \downarrow$ (probability decreases) with ϵ grows & more it

(increase in ϵ)

Proof of theorem :-

$$\text{Def. } C_{n,s} = \sqrt{\frac{2 \ln(n)}{s}}$$

{ In : if let x is arbitrary variable with scalar value

$$P[x \in [x - C_{n,s}, x + C_{n,s}]] \geq 1 - \frac{1}{n}$$

20th Aug '25

$$T_i(n) = 1 + \sum_{m=k+1}^n \{I_m : i\}$$

of top binance minimum in between had some probability

$$\leq m + \sum_{m=k+1}^n \{I_m : i, T_i(m-1) \geq d\}$$

\Rightarrow $\leq m + \sum_{m=k+1}^n \{I_m : i, T_i(m-1) \geq d\}$ showed correctness of analysis
and for i transition see below following set of transitions will

$$\leq d + \sum_{m=k+1}^n \{Q_{m-1}(a^*) + C_{m-1}, T_{a^*}(m-1) \leq$$

performed see above val. \rightarrow quantified to make up to

$$\text{last step } \{Q_{m-1}(a^*) + C_{m-1}, T_{a^*}(m-1) \leq d\}$$

$$P(S_n \geq \mu + E) \leq e^{-2E^2 n}$$

$$P(S_n \leq \mu - E) \leq e^{-2E^2 n}$$

at each & every timestamp, \rightarrow $\mu - E = 0.2$ &

whatever the minimum estimated value of a^*
in particular cycle is less than forward

and estimated value, we showed with probability
less than ϵ

$$\Rightarrow d + \sum_{m=k+1}^n \left\{ \min_{0 \leq s \leq m} (Q_s(a^*) + C_{m-1}, T_{a^*}(s)) \leq d + \epsilon \right\}$$

$$\text{lower bound } \leq \max_{1 \leq s \leq n} (Q_s(i) + C_{m-1}, T_i(s))$$

probability and this satisfies

satisfy now

in each & every step there are a^* values

making upper bound for a particular count

$$\Rightarrow d + \sum_{m=1}^{\infty} \sum_{s=1}^{m-1} \sum_{i=1}^{m-1} \left\{ Q_s(a^*) + C_{m-1}, T_{a^*}(s) \leq Q_s(i) + C_{m-1}, T_i(s) \right\}$$

whatever the probabilities we are getting

$$E[T_i(*)] \leq \left[\frac{8dnCm}{\Delta i^2} \right] + \sum \sum \sum 2m^{-4}$$

3 conditions:

→ under counting / estimated $\hat{\mu}_{\text{arm}}$ $\leq \mu_{\text{optimal}}$ \Rightarrow we can make a mistake.

$$\left[\frac{\text{counts}}{\text{timesteps}} \right] = 1 \quad \text{③}$$

D) underestimated

optimal R then

we can make a mistake.



$$\hat{\mu} = 1A - (1) * p = (\bar{\mu}) * p$$

mistake $\hat{\mu} = 2A$

D) overestimated → can't stop so make a mistake

3) true value of optimal & suboptimal arm's mean is almost neares even in that case we can make a mistake.

$$\hat{\mu} = 1A - (1) * p = (\bar{\mu}) * p$$

→ Underestimating a^*

$$① Q(a^*) \leq q_{\pi}(a^*) - C_m, T_{a^*}(s_i)$$

$$② Q(s_i) \geq q_{\pi}(s_i) + C_m, T_{\pi}(s_i)$$

$$③ q_{\pi}(a^*) \leq q_{\pi}(s_i) + q C_m, T_{\pi}(s_i)$$

concentration bounds.

↳ uncertainty is more b/w optimal & suboptimal so, the reward for it is not known clearly & tend to take a wrong choice.

$$① \boxed{P(Q_s(a^*) \leq q_{\pi}(a^*) - C_m, T_{a^*}(s))} \xrightarrow{\text{confidence interval}} \boxed{\leq m^{-4}} \xrightarrow{\text{eq } n}$$

a^* value at Q_s timestamp.

solving this eqⁿ with

$$P(S_n \geq \mu + \epsilon) \leq e^{-\epsilon^2 n}$$

$$\text{we get } e^{-4 \ln(m) \times p} = e^{-4 \ln m}$$

$$\text{② } \mathbb{P}(\theta_{\text{SI}}(i) \geq q_*(i) + C_m, T_i(s_i)) \leq m^{-4}$$

$$d = \left[\frac{8dn(m)}{\Delta i^2} \right]$$

$$\boxed{q_*(a^*) - q_*(i) - 2dmT_i(s_i)}$$

solve
this eqn with

we get this
value as

$$C_{n,s} = \sqrt{\frac{2dnm}{s}}$$

$$\boxed{q_*(a^*) - q_*(i) - \Delta i = 0.}$$

suboptimal regret

$$\text{Regret}_n = E \left[\sum_i T_i(a) \right] \Delta i$$

$$(12) R_m + (13) \rho \leq (14) \rho$$

$$(12) R_m + (13) \rho \leq (14) \rho$$

$$(12) R_m + (13) \rho \leq (14) \rho$$

previous methods with Action/Value methods.

both worked well w.r.t. bias

works better w.r.t. variance

21st Aug'18

Gradient - Bandit Algorithm:
↳ concerns with every preference action.

$H_t(a)$ → preference function at particular action at time t

initial preference $\rightarrow 0$ for all actions
~~zero~~ ~~so that nothing goes on~~

$\pi_t(a) \rightarrow$ softmax function



$p_t(A_t = a)$

$(H_t(a) = 0) \rightarrow$ initial preference of all actions
for all actions.

gradient ascent / descent preference can be changed
via these steps:

→ for one action - inc de preference

→ for all other actions → dec preference based on reward

baseline - \bar{R}_t (mean) $\neq 0$ in this case

$x_t \rightarrow$ to introduce baseline.

Bandit Problem:

MDPs → useful to represent full RL problems
in a more formal way.

25th Aug '25

States: State / environment in Bandit problem → environment
 Action on same environment.

It's like here → different states which can be merged on the basis of prev time steps and action that we took.

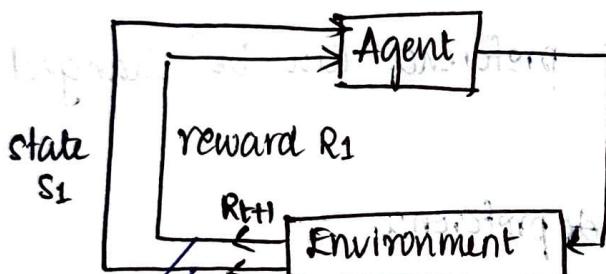
→ Delayed reward

→ Return

→ Value function (wrt each state)

→ Evaluation & feedback (on the basis of each & every state)

full RL problem



Action at time stamp t

state
s_t

reward R_t

R_{t+1}

Environment

get reward + state at next time stamp.

Finite MDP:

$$P(s_t = s' | s_{t-1} = s, A_{t-1} = a) = P(s', r | s, a)$$

what does this function describe?

→ info about 'complete' RL environment!

Agent - Environment Interface:

$$\sum_{s \in S, r \in R} \sum_{s' \in S'} P(s', r | s, a) = 1$$

we can sum probabilities for every possible value of state and reward

* A state having markov property \rightarrow this state has all information about all past agent environment information that matter

for every state - we define some set of actions.

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

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

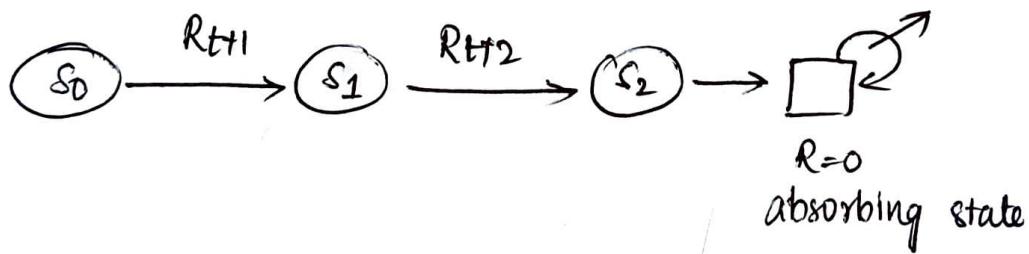
$$0 \leq \gamma \leq 1$$

$$\sum_{k=0}^{\infty} \gamma^k R_{k+t+1}$$

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

$$= R_{t+1} + \gamma (R_{t+2} + \gamma R_{t+3} + \dots)$$

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



POLICY: mapping of states to actions.

$\pi(a|s)$

$$V_{\pi}(s) = E [G_t | s_t = s]$$



state value function.

return $\{ q_{\pi}|a,s = E [G_t | A_t = a, s_t = s] \}$