

1) imagine an investor trying to optimize your trading strategy for ~~your~~ four different stocks labeled A, B, C, D. Each stock has its own unique potential to profit, which is unknown for you. To maximize your return over a series of 100 trades, you decide to implement an ϵ -greedy strategy with ϵ being 0.1. the actual return from each stock follow

Stock A: 70% chance of +1 return,
30% chance of 0.

Stock B: 50% chance of +2 return
50% chance of 0.

Stock C: 10% chance of +5 return
90% chance of 0.

Stock D.: guaranteed return of +0.5

Ans (a) show how do you model this as a
Reinforcement Learning problem

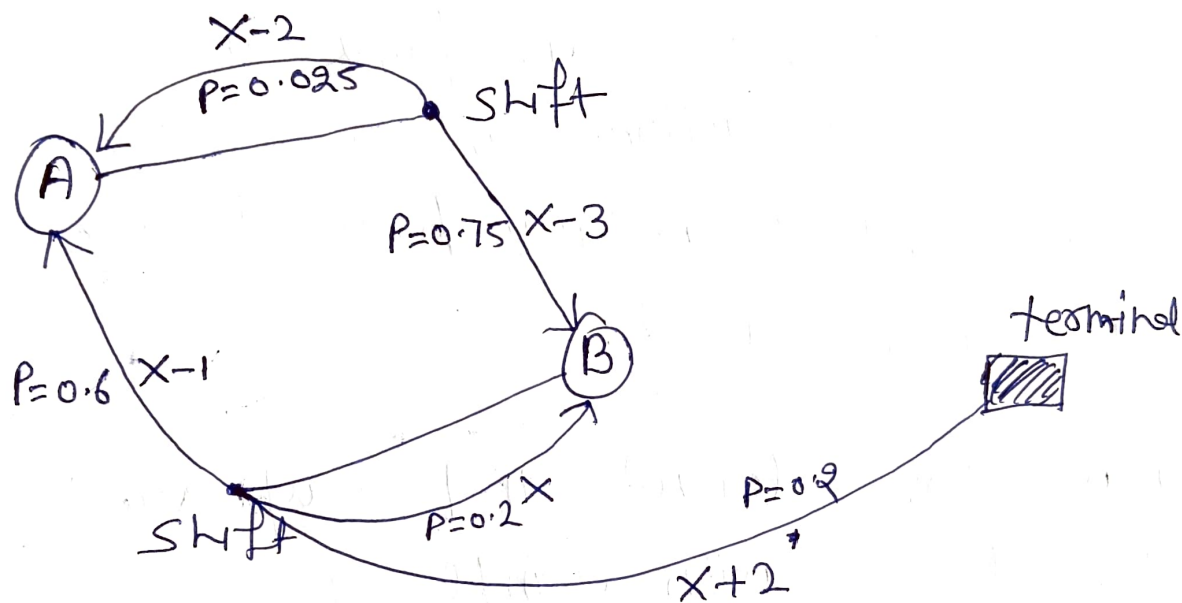
Ans) b) An investor intends to buy 100 times
(each time buying one share of one
stock)

The strategies the investor could

- (i) follow ϵ -greedy for initial 25 trading
and only exploit the information for
next 75 purchase
- (ii) follow ϵ greedy for initial 75 trading
only exploit 25 next purchases
- (iii) follow ϵ greedy for all 100 purchase,
support the investor with your analysis.
show all step, tabulate and to all the
option and write conclusion

(c) What MDP, POMDP and CMDP? What are
they? Suggest one RL technique that is
used to solve problem study with them
It is adequate if you write just one line
for each

Q-2: Consider the MDP given below containing 2 states A & B with action shift that may result in A, B or terminal state. The Reward obtained are indicated along the edges in the figure ($x-2$, $x-3$, $x-1$, x , $x+2$). Let x to be 6. The transition probabilities are also given along the edges. Discount factor = 0.4.



- (a) Evaluate the given deterministic policy where the shift always executes the higher probability action. Improve it upto 1 iteration, use dynamic programming solution to MDP.
- (b) With value iteration of dynamic programming determine the value of states A and B. Let the value of A & B initialized to 1, show 1 iteration

Q=3(a) what are the most important issues when you have to learn the value function using a first visit Monte Carlo using for a deterministic policy. Explain also provide possible solution.

Q=3(b) Explain any 3 most significant action selection strategies used in RL and mention how each selection method balances exploration and exploitation. provide your ans in table

Q=3(c) if we utilize policy gradient method to address a Reinforcement learning model and find that policy is provide not optimal, what could be the possible explanation for this. state method
select 3 reason

Q24

(a) Why AlphaGo use a ^{Separate} ~~Separate~~ policy network and a separate value network.

(b) How does the MCTS ensure an action with the highest value is found in real time? If the best action can be selected only by MCTS, why is any prior learning of $Q(S, a)$ required.

(c) We have learned that supervised learning that learns with samples from a given distribution does not capture the online ~~value~~ nature of interactions as required for reinforcement learning quite well.

(i) Why does AlphaGo use supervised learning to learn the initial policy (and even further)
ii) In what way the shortcomings of supervised learning are mitigated in AlphaGo?

(d) How does DQN handle the challenge reflected in the part (c) of this question

5 a) Consider the following ways of organizing Reinforcement learning techniques.

- (i) Model based vs Model free
- (ii) value based vs policy based
- (iii) on policy vs off policy.

Write a statement or two on each of the points (for both category) explaining the kind of problem these RL technique are suited to. provide your response in ~~neat~~ neatly organized table.

5(b) Consider the learning scenario. A human expert is presented with two trajectories taken by two drivers in the highway stretch. The human expert marks which of the trajectory better. The agent learns this expertise (to decide a better trajectory by giving two unseen trajectory) observing the expert decision from many such example. Explain how you precisely

model this as an appropriate RL problem.
show all the elements of your model.
Making necessary assumptions