

HOMWORK - 1

1. CONCEPT CHECK

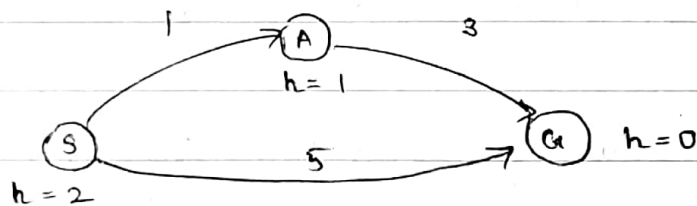
1. Data structures for collection of fringe / frontier nodes

BFS \rightarrow FIFO

DFS \rightarrow LIFO

A* search \rightarrow Priority Queue

2.



admissible and consistent heuristic

$$h(S) = 2 \leq \text{actual cost (S to G)}$$

$$2 \leq 4$$

$$h(A) = 1 \leq \text{actual cost (A to G)}$$

$$1 \leq 3$$

$$h(S) - h(G) < \text{cost (S to G)}$$

$$2 - 0 \leq 4$$

$$2 \leq 4$$

$$h(S) - h(A) \leq \text{cost (S to A)}$$

$$2 - 1 \leq 1$$

$$1 \leq 1$$

$$h(A) - h(G) \leq \text{cost}(A \text{ to } G)$$

$$1 - 0 \leq 3$$

$$1 \leq 3$$

$$0 \leq h(n) \leq h^*(n)$$

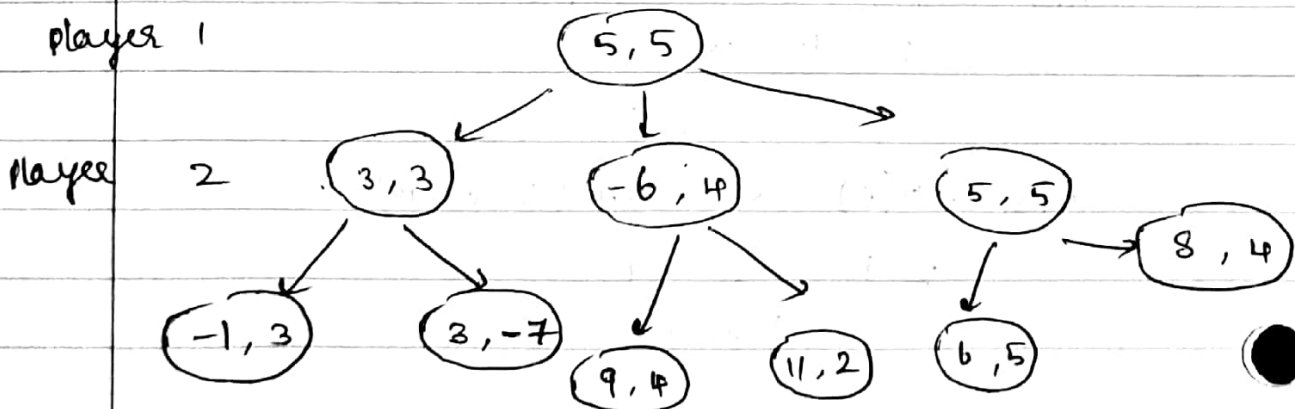
$$0 \leq h(s) \leq h^*(s)$$

$$0 \leq 2 \leq 4$$

$$0 \leq h(A) \leq h^*(A)$$

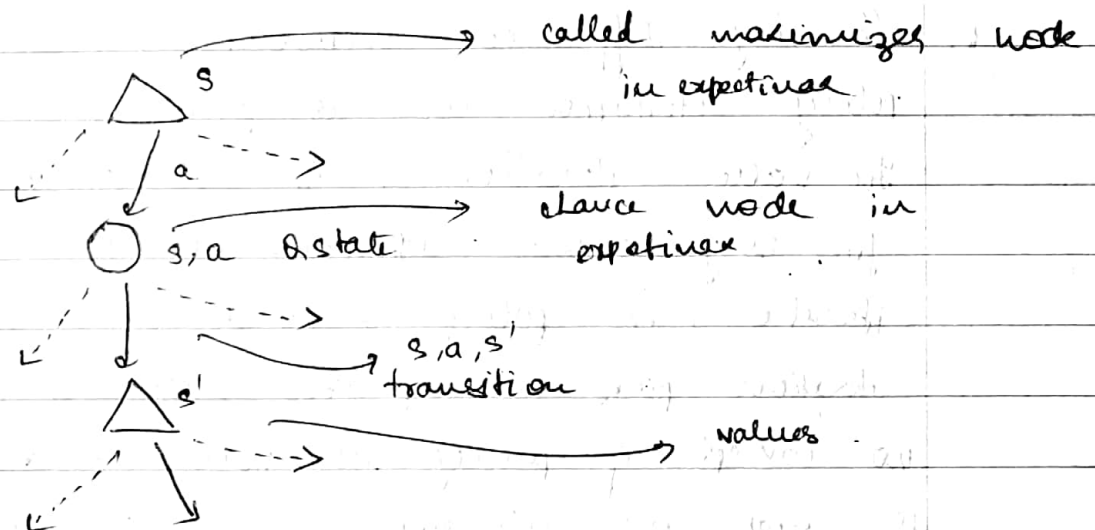
$$0 \leq 1 \leq 3$$

3. In a non-zero game setting, we can't make use of alpha-beta pruning. This is because the players trying to maximize the values are independent. Which implies there is no more the case where one player tries to let down other player.



In the above tree, there is no minimizer node and only 2 players who play to maximize their values are present.

4. The relationship between expectimax algorithm and Q -states working in MDPs are as follows:



A maximizer node is the state (s) in MDP. The when we take an action (a) we land in a state. We call this as a state node in expectimax. A random transition will happen based on probability values and the resulting state s' .

5. The assumption that our decision processes are Markov based can be attributed to the inherent goal of MDP which is to formulate a policy that specifies an action to be taken with respect to current state. The Markov

property is necessary since the policy aims to increase the cumulative rewards as a function of current state. Hence, only the current state plays a role in deciding the outcome which is directly what we expect of a policy.

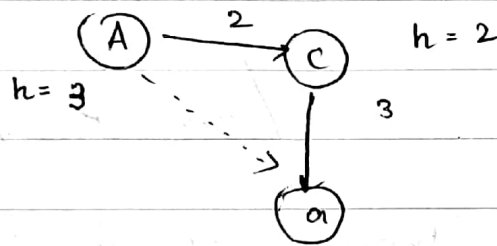
1. The key differences between value and policy iteration are as follows:
- In value iteration, the values are updated for every round whereas in the policy iteration, the policy is fixed and multiple iterations pass in updating values. There is no concept of policy in value iteration, only the next one all values is the action to be taken. In policy iteration, all actions are taken as per the policy.

3. FUN WITH PROOFS

1. ADMISSIBILITY : $h(A) \leq \text{actual cost from } A \text{ to } G$
CONSISTENCY : $h(A) - h(C) \leq \text{cost}(A \rightarrow C)$

\therefore All consistent heuristics will be admissible

This can be showed as follows :-



consistency is

$$h(A) - h(C) \leq \text{cost}(A \text{ to } C)$$

$$\text{cost}(A \text{ to } C) = \text{cost}(A \text{ to } G) - \text{cost}(C \text{ to } G)$$

$$\text{cost}(A \text{ to } C) = h(A) - h(G)$$

if they are same.

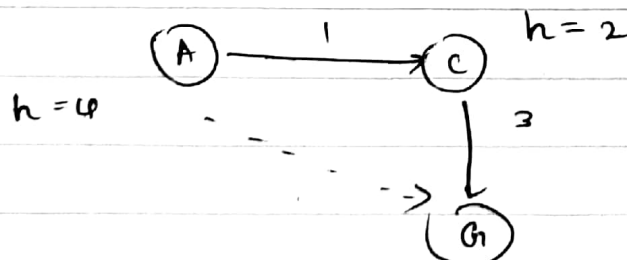
$$h(A) - h(C) \leq \text{cost}(A \text{ to } G) - \text{cost}(C \text{ to } G)$$

$$h(A) \leq \text{cost}(A \text{ to } G)$$

$$h(C) \leq \text{cost}(C \text{ to } G)$$

Automatically they turn out to be admissible

However, not all admissible heuristic are consistent.



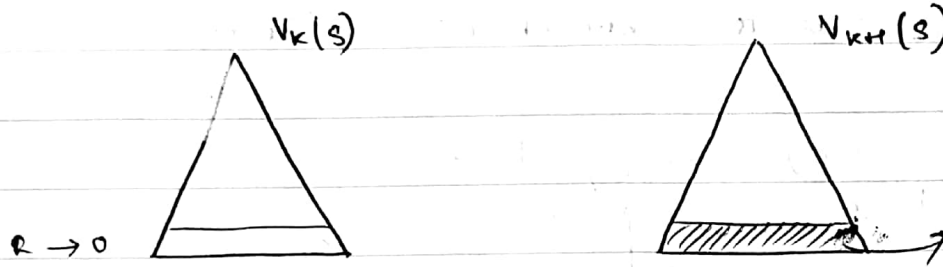
$$h(A) - h(C) \leq \text{cost}(A \text{ to } C)$$

$$4 - 2 \leq 1$$

$$2 \neq 1$$

hence, proved.

2.



The only difference between V_k and V_{k+1} is the last layer which is

$$[V_k \times R_{MAX} \text{ (or) } V_k \times R_{MIN}]$$

$$V_{k+1}(s) = \max_a \sum_{s'} \tau(s, a, s') [R(s, a, s') + \gamma V_k(s')]$$

γ^k because for each step the discount factor is raised accordingly

When the discount is less than 1, the resulting values in the last layer will settle since as k approaches infinity the value will approach zero or constant.