

CS747 - Foundations of Intelligent and Learning Agents

CS747 Autumn 2023
Chinmay Makarand Pimpalkhare
200100115

Assignment 2 | [Markov Decision Planning Methods for Half-field offense Football](#)

1 Task 1

Default algorithm used is Value Iteration, due to fastest convergence even on the large MDP

1.1 Value Iteration

1.1.1 Bellman Optimality Operator

Implementation:

```
1 def Bellman_optimality_operation(num_states, num_actions, mdp_type, is_mdp_episodic, \
2     these_are_episodic, transition_matrix, reward_matrix, \
3     gamma, value_function
4 ):
5     #Steps: For all states, max over actions (sum of (T*(R + gamma*V)))
6     new_value_function = np.zeros(num_states)
7     optimal_action = np.zeros(num_states)
8
9     for i in range(num_states):
10         if these_are_episodic[i] == 0:
11             new_value_function[i] = np.max(np.sum(np.multiply(transition_matrix[i], \
12                 reward_matrix[i] +\
13                     gamma*value_function), axis = 1))
14             optimal_action[i] = np.argmax(np.sum(np.multiply(transition_matrix[i], \
15                 reward_matrix[i] +\
16                     gamma*value_function), axis = 1))
17
18     return new_value_function, optimal_action
```

Details:

- The new value function V has been initialized to 0 before applying the operator to it
- The value functions of only non-terminal states are updated. Thus, it is guaranteed that terminal states will always have a value function which is identically 0

1.1.2 Value Iteration

Implementation

```
1 def value_iteration(num_states, num_actions, mdp_type, is_mdp_episodic, \
2     these_are_episodic, transition_matrix, reward_matrix, \
3     gamma):
4     epsilon = 1e-12
5     V_old = np.zeros(num_states)
6     V_new, pi = Bellman_optimality_operation(num_states, num_actions, mdp_type,
7         is_mdp_episodic, \
```

```

7     these_are_episodic, transition_matrix, reward_matrix, \
8     gamma, V_old)
9     V_copy = V_old
10    while np.linalg.norm(V_new - V_copy, 1) >= epsilon:
11        V_copy = V_old
12        V_new, pi = Bellman_optimality_operation(num_states, num_actions, mdp_type,
13        is_mdp_episodic, \
14        these_are_episodic, transition_matrix, reward_matrix, \
15        gamma, V_old)
16        V_old = V_new
17    for i in range(num_states):
18        print(V_new[i], int(pi[i]))

```

Details:

- The initial value function at the start of solving the MDP, V_0 has been initialized to $\mathbf{0}$
- The Bellman optimality operator is then repeatedly applied to the value function to get the updated value function.
- The termination criteria used is that the \mathcal{L}_1 norm of $V^t - V^{t+1}$ should be less than or equal to $\varepsilon = 10^{-12}$

1.2 Linear Programming

Implementation

```

1  def linear_programming(num_states, num_actions, mdp_type, is_mdp_episodic, \
2  these_are_episodic, transition_matrix, reward_matrix, \
3  gamma):
4  prob = LpProblem("Linear_Programming_Value", LpMaximize)
5  set_S = range(0, num_states)
6  set_A = range(0, num_actions)
7  val_func = LpVariable.dicts("Value_Function", set_S, cat='Continuous')
8
9  #objective function for the primal
10 prob += lpSum([-val_func[i] for i in set_S])
11
12 #nk constraints
13 for i in set_S:
14     for j in set_A:
15         prob += val_func[i] >= lpSum([transition_matrix[i][j][k]*\
16         (reward_matrix[i][j][k] + \
17         gamma*val_func[k]) for k in set_S])
18
19
20 #solve
21 prob.solve(PULP_CBC_CMD(msg=0))
22 V_star = np.array([val_func[i].varValue for i in set_S])
23 Q_values = np.zeros((num_states, num_actions))
24 for i in set_S:
25     for j in set_A:
26         Q_values[i][j] = np.sum(np.multiply(transition_matrix[i][j], \
27         reward_matrix[i][j] + \
28         gamma*V_star))
29
30 pi_star = np.argmax(Q_values, axis = 1)
31 for i in set_S:
32     print(V_star[i], int(pi_star[i]))

```

Details:

- Objective function used is (as given in slides)

$$\min - \left(\sum_{s \in S} V(s) \right)$$

- nk Constraints used are also as given in the slides

- The primal problem is solved and to compute the optimal policy, we compute Q-values and pick the action with maximum Q-value for a given state
- Ties are broken internally by np.argmax

1.3 Howard's Policy Iteration

Implementation

```

1  def howard_policy_iteration(num_states, num_actions, mdp_type, is_mdp_episodic, \
2      these_are_episodic, transition_matrix, reward_matrix, \
3      gamma):
4
5      #initialize arrays for improvable states and improvable actions
6      #corresponding to the improvable states
7      improvable_states = np.ones(num_states)
8      improving_actions = np.zeros((num_states, num_actions))
9      pi_old = np.random.randint(num_actions, size = num_states)
10     while(np.sum(improvable_states) > 0):
11         V_old = policy_evaluation_no_print(num_states, num_actions, mdp_type,
12             is_mdp_episodic, \
13             these_are_episodic, transition_matrix, reward_matrix, \
14             gamma, pi_old)
15         Q_pi = compute_action_values(num_states, num_actions, mdp_type, is_mdp_episodic, \
16             these_are_episodic, transition_matrix, reward_matrix, \
17             gamma, V_old)
18         improvable_states.fill(0)
19         improving_actions.fill(0)
20         for i in range(0,num_states):
21             for j in range(0,num_actions):
22                 if Q_pi[i][j] > V_old[i] and j != pi_old[i] :
23                     improvable_states[i] = 1
24                     improving_actions[i][j] = 1
25         for i in range(0,num_states):
26             if improvable_states[i] == 1:
27                 pi_old[i] = int(np.random.choice(num_actions, \
28                     1, p = improving_actions[i]/np.sum(improving_actions[i])))
29
30     for i in range(0, num_states):
31         print(V_old[i],int(pi_old[i]))

```

Details

- The arrays pertaining to improvable states and actions are initialized by 1 and 0 respectively
- The initial policy π_0 is randomly chosen
- Termination occurs when $IS(s) = \phi \ \forall s \in S$
- Policy evaluation is used for computing the value function of the given policy. It uses PuLP in the background to solve a system of n **linear equations**. (First method shown in class)
- Inside the loop, all states with improvable actions are marked by comparing action values with value functions. The improvable actions are stored
- **For all improvable states**, we switch randomly to a new action using the np.choice function

1.4 Policy Evaluation

Implementation

```

1  # A modification of policy evaluation which only returns the value functions
2  def policy_evaluation_no_print(num_states, num_actions, mdp_type, is_mdp_episodic, \
3      these_are_episodic, transition_matrix, reward_matrix, \
4      gamma, policy):
5      #We shall be using pulp
6      #We can define a dummy objective function

```

```

7  #We have the following 2 constraints for each state:
8  #   V(s) >= \sum_{s'} T*(R + V(s'))
9  #   V(s) <= \sum_{s'} T*(R + V(s'))
10 prob = LpProblem("Policy_Evaluation")
11 set_S = range(0, num_states)
12 val_func = LpVariable.dicts("Value_Function", set_S, cat='Continuous')
13 for i in set_S:
14     prob += val_func[i] == lpSum([ transition_matrix[i][int(policy[i])][j]*\
15                                     (reward_matrix[i][int(policy[i])][j] + \
16                                     gamma*val_func[j]) for j in set_S])
17 prob.solve(PULP_CBC_CMD(msg=0))
18 V_pi = np.array([val_func[i].varValue for i in set_S])
19 return V_pi

```

Details

- Define n equations corresponding to the n states of the MDP. Solves them using PuLP to get the value function for a given **deterministic** policy

1.5 Summary of Algorithms

- On given test cases for task 1, HPI and LP were faster than VI. However, they were very slow for the football problem. Hence, I found that on smaller MDPs it was convenient to use HPI and LP
- VI was very slow on test case 4, possibly because the value functions were large in magnitude

2 Task 2: 2v1 Football

Default algorithm used is Value Iteration, due to best performance.

2.1 MDP Formulation

States There are in all 8194 states. 8192 states are non-terminal and correspond to the $16 \times 16 \times 16 \times 2$ string configurations possible as given in the policy files. The last 2 states are terminal. The 8193th state corresponds to termination due to losing the game with a reward 1, while the 8194th state corresponds to scoring a goal and the game ending with a reward equal to 1.

```

1  def state_encoder(state_string):
2  if state_string[0] == '0':
3      B1_state = int(state_string[1])
4  elif state_string[0] == '1':
5      B1_state = int(state_string[0:2])
6  if state_string[2] == '0':
7      B2_state = int(state_string[3])
8  elif state_string[2] == '1':
9      B2_state = int(state_string[2:4])
10 if state_string[4] == '0':
11     R_state = int(state_string[5])
12 elif state_string[4] == '1':
13     R_state = int(state_string[4:6])
14 poss = int(state_string[-1])
15 state_encoding = int((B1_state - 1)*512 + (B2_state - 1)*32 + (R_state - 1)*2 + poss - 1)
16 return state_encoding

```

State encoding is done for convenience as shown in the code above. It will be useful for the actions.

Actions We use the actions as given in the question. For checking if a given action is legitimate (meaning that the players do not go out of bounds), we define left, right, up and down operators which are functions of the string description of a state and a number which tells us which player has moved)

```

1 def left_operator(state_string, move_whom): #Move_whom is 1 for B1, 2 for B2 and 3 for R
2     B1 = state_string[0:2]
3     B2 = state_string[2:4]
4     R = state_string[4:6]
5     poss = state_string[-1]
6     if move_whom == 1:
7         if B1 == '02':
8             new_state = '01' + B2 + R + poss
9             return new_state
10        elif B1 == '03':
11            new_state = '02' + B2 + R + poss
12            return new_state
13        elif B1 == '04':
14            new_state = '03' + B2 + R + poss
15            return new_state

```

The first few lines of the "left_operator" function have been given. The logic is the same for other 3 movements.

Tackling, Shooting and Passing Appropriate conditions have been implemented in the code to check for all these operations (refer to code)

Transitions The transitions from state s to s' following action a occur with probability

$$\begin{aligned}
 T(s, a, s') &= \sum_{a_0 \in L, R, U, D} T(s^*, a, s') \mathbb{P}(a_0) \\
 &= \sum_{a_0 \in L, R, U, D} T(s^*, a, s') \pi_0
 \end{aligned}$$

where we are conditioning on the action that the opponent takes, and thus the policy of the opponent

Rewards If a goal is scored, a reward of 1 is obtained. All other transitions have reward equal to 0.

Discount Factor We assume a discount factor equal to 1, since this is an episodic MDP

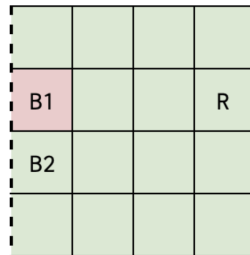


Fig 11: Starting state

Figure 1: Starting state for the MDP as per the evaluation metrics

Computing Expected Number of Goals

- First, we must note that the expected number of goals starting from a given state is nothing but the value function of that state. This is because the reward is 1 only when a goal is scored, thus if we examine the equation

$$\mathbb{E}[\text{number of goals scored starting at } s] = \mathbb{P}(\text{event that we end in a goal-scoring terminal state})$$

- Further,

$$\mathbb{E}[\text{number of goals scored starting at } s] = \mathbb{E}[(R_0 + R_1 + R_2 + \dots + R_T | s_0 = s)]$$

where $R_0, R_1, R_2, \dots, R_T$ are the rewards obtained (goals scored) at timesteps given by the subscripts

- Note that as per our MDP formulation, only scoring a goal gives us a goal. Thus $\forall t < T, R_t = 0$
- Hence, we have

$$\mathbb{E}[\text{number of goals scored starting at } s] = \mathbb{E}[(R_T | s_0 = s)]$$

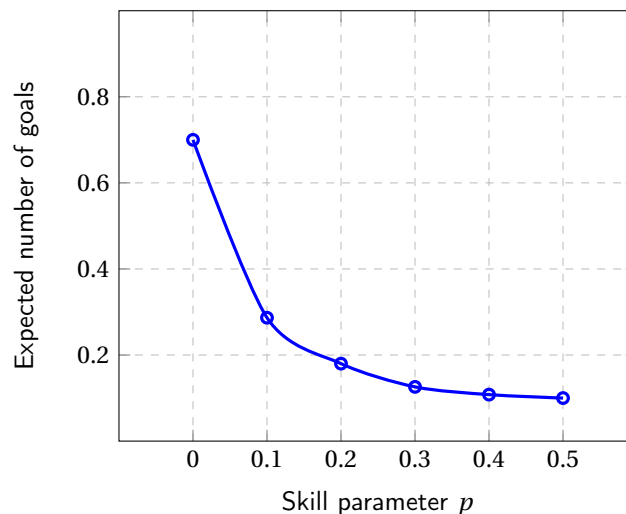
- But note that $\mathbb{E}[(R_0 + R_1 + R_2 + \dots + R_T | s_0 = s)]$ is nothing but the value function of state s given policy π assuming $\gamma = 1$
- Thus we have

$$\mathbb{E}[\text{number of goals scored starting at } s] = V^\pi(s)$$

- Thus, the task reduces to computing $V^\pi(s)$

2.2 Graph 1: Variation of expected number of goals with p , keeping q constant against greedy defense (policy 1)

Variation when $q = 0.7$ is constant and opponent follows greedy defense



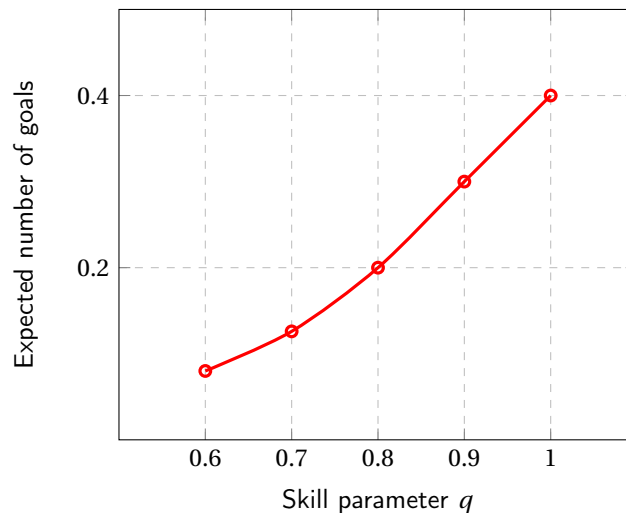
Observations

- From the graph, we can see that $V^\pi(s)$ reduces as p increases. This can be seen from the following statement:
 - Probability of losing the game by unsuccessful movement is proportional to p , with the constant of proportionality depending on who has possession (with ball or without ball).
- If $p = 0$, then almost all movements (except the tackles) are successful. This case has a very high value function
- The policy which the opponent is following is greedy defense. Thus, he is trying to move towards the player with the ball. This increases the chances of there being tackles. Hence, parameter p is of great importance as it mainly controls movement.

- The rate of change of reduction in value function gets smaller as we increase p . As p increases, movements start becoming more and more unsuccessful so we are more likely to lose possession.
- If we analyse the **optimal** actions being taken at this state, we see the following things
 - The optimal actions are 1, 1, 1, 5, 5, 9 corresponding to the six p values as given.
 - Thus, if p is very low, the player with the ball B_1 moves. The strategy might be to allow the defender to come towards him and make a pass to the adjacent player, since now the defender has been pulled out from the goal, which increases the probability of a successful strike to goal. Also, the player B_2 can sneak into the goal.
 - Also, R is in front of B_1 and there is high possibility of tackle after some iterations of the game. At $p = 0.3$ onwards, the optimal action changes and B_2 starts to move. This is because $2p$ has become very large so the player with the ball moving becomes very risky.
 - As p increases to 0.5, movement of B_1 will always lead to failure and even movement of B_2 will fail with probability 0.5, hence the best choice is to shoot. For this case, most other states also have the optimal action to be shoot

2.3 Graph 2: Variation of expected number of goals with q , keeping p constant against greedy defense (policy 1)

Variation when $p = 0.3$ is constant and opponent follows greedy defense



Observations

- From the graph, we can see that $V^\pi(s)$ increases as q increases. This can be seen from the following statements:
 - The probability is a function of the distance from the teammate for passing the ball. If player B1 is on (x_1, y_1) and B2 is on (x_2, y_2) then the probability of a successful pass is $q - 0.1 * \max(|x_1 - x_2|, |y_1 - y_2|)$
 - For shooting towards the goal, the probability is a function of the distance of the x coordinate of the shooter. If B1 has the ball and is on (x_1, y_1) the probability of a goal is $q - 0.2 * (3 - x_1)$.
- Thus q directly affects the chance that a pass or a goal is successful and with increase in q , the success rate increases.
- Again the intuition is that movement is risky in case of greedy defense since the defender is directly moving towards the player who has the ball, and passing is more likely to make the defender switch directions. Hence, success rate of passing is important.

- We can also analyse the **optimal actions** taken for each value of q
 - The optimal actions taken at the state are 1, 5, 9, 9, 9 for the five values of q given
 - This is because there is always a trade-off between movement and passing/shooting. Passing and shooting is safer as compared to movement at larger values of q and hence as q increases the tendency to shoot increases.
 - The rate of increase in expected number of goals increases as q increases. This is because the optimal action gets fixed after a certain threshold of q and thus we can relate the expected number of goals directly to the shooting accuracy. In the initial phase of the curve, since the optimal action is movement, p is the defining parameter.