

AI6101 Reinforcement Learning Assignment

Goh Jun Hao G21202327

```
In [1]: from typing import List, Tuple
import copy
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
%matplotlib inline
```

The following class defines the grid world environment. The grid world looks like:

	0	_1_	_2_	_3_	_4_	_5_	_6_	_7_	_8_	_9_	_10_	_11_	_12_	_13_
0							x	x						
1							x	x						
2				x			x	x					x	
3				x			x					x	x	
4		B		x								x	x	G
5	A			x								x	x	

```
In [2]: AGENT = 'A'
BOX = 'B'
GOAL = 'G'
DANGER = 'x'
GRID = '_'

class CliffBoxGridWorld:
    """
    Cliff Box Pushing Grid World.
    """
    action_space = [1, 2, 3, 4]
    forces = {
        1: np.array([-1, 0]),
        2: np.array([1, 0]),
        3: np.array([0, -1]),
        4: np.array([0, 1]),
    }
    world_width = 14
    world_height = 6
    goal_pos = np.array([4, 13])
    init_agent_pos = np.array([5, 0])
    init_box_pos = np.array([4, 1])
    danger_region = [
        [(2, 3), (5, 3)],
        [(0, 6), (3, 6)],
        [(0, 7), (2, 7)],
        [(3, 11), (5, 11)],
        [(2, 12), (5, 12)],
    ]

    def __init__(self,
                  episode_length=100,
                  render=False,
                  ):
        """
```



```

    # Calculate the rewards
    done = self.timesteps == self.episode_length - 1
    # the distance between agents and box
    dist = np.sum(np.abs(self.agent_pos - self.box_pos))
    reward = -1 # -1 for each step
    reward -= dist
    # if agents or box is off the cliff
    if self._check_off_cliff(self.agent_pos) or self._check_off_cliff(self.box_pos):
        reward += -1000
        done = True

    if all(self.box_pos == self.goal_pos):
        reward += 1000
        done = True

    reward -= np.sum(np.abs(self.box_pos - self.goal_pos))

    if self.render:
        self._update_render()

    self.timesteps += 1
    info = {}

    return state, reward, done, info

def print_world(self):
    """
    Render the world in the command line.
    """
    if len(self.action_history) > 0:
        print(f'Action: {self.action_history[-1]}')
    print(self.world)

def _check_pos_boundary(self, pos, box_hard_boundary: bool = False):
    """
    Move the given position within the world bound.
    """
    if pos[0] < 0:
        pos[0] = 0
    if pos[0] >= self.world_height:
        pos[0] = self.world_height - 1
    if pos[1] < 0:
        pos[1] = 0
    if pos[1] >= self.world_width:
        pos[1] = self.world_width - 1

    if box_hard_boundary:
        if pos[0] == 0:
            pos[0] += 1
        elif pos[0] == self.world_height - 1:
            pos[0] = self.world_height - 2
        if pos[1] == 0:
            pos[1] += 1

    return pos

def _check_off_cliff(self, pos):
    """
    Check if the given position is off cliff.
    """
    for region in self.danger_region:
        A, B = region
        assert A[1] == B[1], "A[1] != B[1]"
        if A[0] <= pos[0] <= B[0] and pos[1] == A[1]:
            return True

```

```

return False

def _update_render(self):
    """
    Update the render information.
    """
    if not all(self.last_agent_pos == self.agent_pos):
        pos = self.last_agent_pos
        if (pos[0] != self.goal_pos[0]) or (pos[1] != self.goal_pos[1]):
            self.world[pos[0], pos[1]] = GRID

    if not all(self.last_box_pos == self.box_pos):
        pos = self.last_box_pos
        if self.world[pos[0], pos[1]].decode('UTF-8') not in {AGENT}:
            self.world[pos[0], pos[1]] = GRID

    if (self.agent_pos[0] != self.goal_pos[0]) or (self.agent_pos[1] != self.goal_pos[1]):
        self.world[self.agent_pos[0], self.agent_pos[1]] = AGENT
    self.world[self.box_pos[0], self.box_pos[1]] = BOX
    self.last_box_pos = copy.deepcopy(self.box_pos)
    self.last_agent_pos = copy.deepcopy(self.agent_pos)

```

Q-learning Agent

Q-learning is a value-based method that uses a Q- table to record the estimated Q-values for various actions in different states. The Q-table is initialized prior to exploring the environment. As the agent interacts with the environment, it updates the Q(s, a) values using the Bellman equation, allowing the agent to continually learn and improve its understanding of the environment. The Q-function continues to be updated until it reaches convergence or the specified number of iterations

```

In [3]: class QAgent:
    def __init__(self, env, num_episodes, epsilon=0.1, alpha=0.1, gamma=0.99):
        self.action_space = env.action_space
        self.q_table = dict() # Store all Q-values in a dictionary
        # Loop through all possible grid spaces, create sub-dictionary for each
        for agent_x in range(env.world_height):
            for agent_y in range(env.world_width):
                for box_x in range(env.world_height):
                    for box_y in range(env.world_width):
                        # Populate sub-dictionary with zero values for possible moves
                        self.q_table[(agent_x, agent_y, box_x, box_y)] = {k: 0 for k in env.action_space}

        self.env = env
        self.num_episodes = num_episodes
        self.epsilon = epsilon
        self.alpha = alpha
        self.gamma = gamma
        self.action_seq = [] # Initialize empty action sequence list
        self.total_action_seq = [] # Initialize empty action sequence list
        self.max_agent_seq = []

    def act(self, state):
        """Returns the (epsilon-greedy) optimal action from Q-Value table."""
        if np.random.uniform(0,1) < self.epsilon:
            action = self.action_space[np.random.randint(0, len(self.action_space))]
        else:
            q_values_of_state = self.q_table[state]
            maxValue = max(q_values_of_state.values())
            action = np.random.choice([k for k, v in q_values_of_state.items() if v == maxValue])

        return action

```

```

def learn(self):
    """Runs Q-learning algorithm to learn optimal policy."""
    rewards_per_episode = []

    for episode in range(self.num_episodes):
        # Initialize variables for storing state and action sequences
        action_seq = []
        state = self.env.reset()
        total_reward = 0
        done = False
        while not done:
            action = self.act(state)
            next_state, reward, done, _ = self.env.step(action)

            # Update Q-Table using Q-Learning algorithm
            old_qvalue = self.q_table[state][action]
            next_max = max(self.q_table[next_state].values())
            new_qvalue = (1 - self.alpha) * old_qvalue + self.alpha * (reward + self
            self.q_table[state][action] = new_qvalue

            # Update total reward and current state
            total_reward += reward
            state = next_state
            action_seq.append(action)
            self.action_seq = action_seq # Append action to the action sequence list

        self.total_action_seq.append(action_seq)
        rewards_per_episode.append(total_reward)
        max_reward = max(rewards_per_episode)
        index_max_reward = rewards_per_episode.index(max_reward)
        self.max_action_seq = self.total_action_seq[index_max_reward]

    #print(rewards_per_episode)
    print('Maximum reward =', max_reward)
    print('Episode of maximum reward achieved =', index_max_reward)
    print('Action sequence of maximum reward =', self.max_action_seq)
    return rewards_per_episode

```

Training Q-learning Agent

```

In [4]: env = CliffBoxGridWorld()
        agent = QAgent(env, epsilon=0.01, alpha=0.1, gamma=0.99, num_episodes=15000)
        rewards = agent.learn()

```

```

Maximum reward = 642
Episode of maximum reward achieved = 6225
Action sequence of maximum reward = [4, 1, 1, 1, 3, 1, 4, 4, 4, 4, 1, 4, 2, 2, 2, 3, 2,
4, 4, 4, 4, 4, 2, 4, 1, 1, 1, 3, 1, 4, 4, 4, 1, 4, 2, 2, 2]

```

Q-Learning Episodes vs Episode Rewards Curve

It can be observed from the plot below that the q-learning agent converge nicely at the maximum reward of 642 by about episode 7000

```

In [5]: # Smooth plot
        weight = 0.99
        last = rewards[0]
        smoothed = []
        for v in rewards:
            smoothed_val = last * weight + (1 - weight) * v
            smoothed.append(smoothed_val)

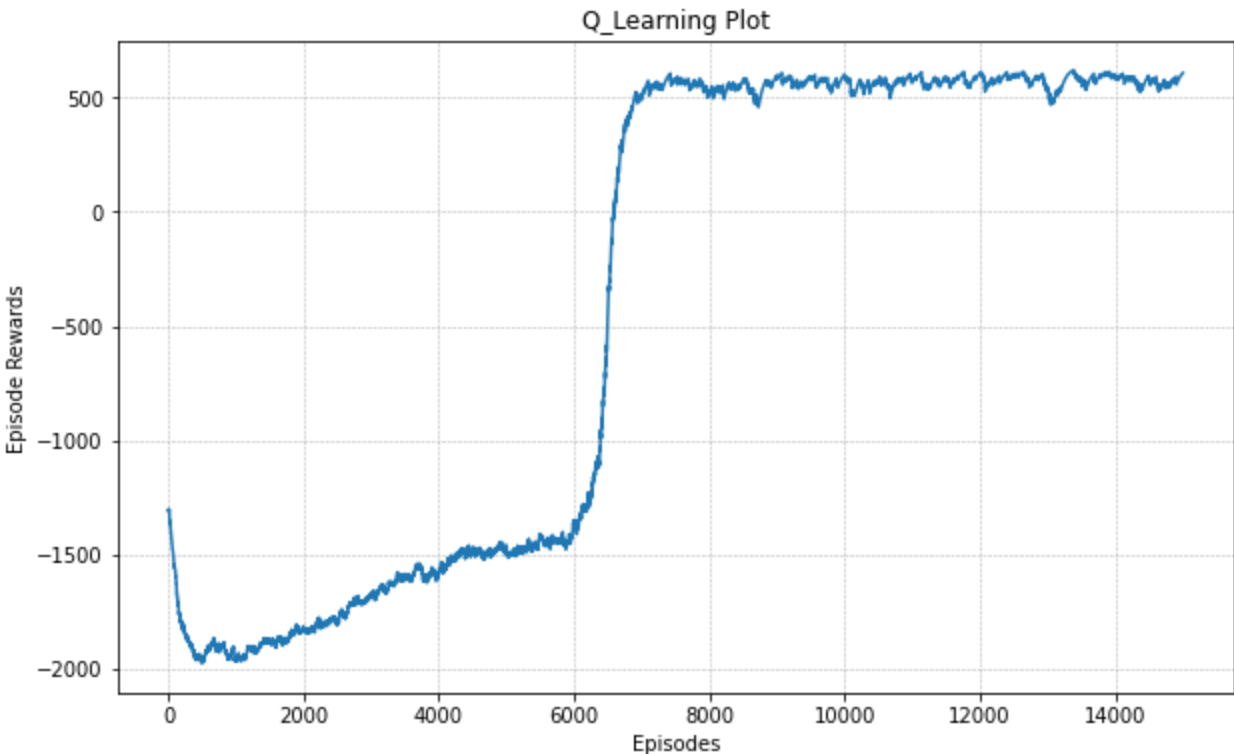
```

```

last = smoothed_val

# Plot the learning curve
plt.figure(figsize=(10, 6))
plt.plot(smoothed)
plt.xlabel("Episodes")
plt.ylabel("Episode Rewards")
plt.title('Q_Learning Plot')
plt.grid(linestyle='--', linewidth=0.5)
plt.show()

```



Q-learning V_table Visualization

Finding the 'V' for the agent

```

In [6]: v_table = {}
        for state_, actions_ in agent.q_table.items():
            v_table['('+','.join(map(str, state_))+')'] = list(actions_.values())
        df = pd.DataFrame.from_dict(v_table).transpose()
        df['state'] = list(df.index)
        df['Agent_Position'] = df['state'].apply(lambda x:x[1:5].rstrip(','))
        df_group = df.groupby('Agent_Position')[[0, 1, 2, 3]].mean()

        df_group = df.groupby('Agent_Position')[[0, 1, 2, 3]].mean().reset_index()
        df_group['V'] = df_group[[0, 1, 2, 3]].mean(axis=1)
        df_group[['agent_x', 'agent_y']] = df_group['Agent_Position'].str.split(',', expand=True)

        # showing the data frame in terms of x and y coordinate of agent, and 'V'
        df_group = df_group[['agent_x', 'agent_y', 'V']]
        df_group

```

```

Out[6]:
   agent_x  agent_y    V
0         0         0 -38.931412
1         0         1 -39.317577
2         0        10 -30.111324

```

3	0	11	-26.464634
4	0	12	-20.204145
...
79	5	5	-37.677345
80	5	6	-37.081845
81	5	7	-35.447405
82	5	8	-34.693796
83	5	9	-28.446331

84 rows × 3 columns

Presentation of V_table in the template of the Cliff Box Grid World

It can be observed from the grid view of the 'V' that the obstacles 'V' are all zero, which by intuition make sense.

```
In [16]: # create a 6x14 array of NaN values
grid = np.empty((6, 14))
grid[:] = np.nan

# fill in the grid with the V values from the dataframe
for _, row in df_group.iterrows():
    x = int(row['agent_x'])
    y = int(row['agent_y'])
    grid[x, y] = round(row['V'], 2)

# create a dataframe from the grid
grid_df = pd.DataFrame(grid, columns=range(14), index=range(6))

# set format for V values
grid_df = grid_df.style.format("{:.2f}")

# add borders and show row/column indices
grid_df = grid_df.set_table_styles([
    {"selector": "th", "props": [("border", "1px solid #ccc"), ("text-align", "center")]},
    {"selector": "td", "props": [("border", "1px solid #ccc"), ("text-align", "center")]},
    {"selector": "th.row_heading", "props": [("background-color", "transparent"),]}
])

# display formatted grid of v_table
print('V_table Visualization')
display(grid_df)
```

V_table Visualization

	0	1	2	3	4	5	6	7	8	9	10	11	12	13
0	-38.93	-39.32	-39.80	-37.75	-33.89	-44.51	0.00	0.00	-41.71	-30.93	-30.11	-26.46	-20.20	-16.49
1	-31.63	-30.82	-32.17	-41.92	-36.60	-43.60	0.00	0.00	-38.03	-23.32	-22.88	-22.87	-32.19	-17.33
2	-35.12	-31.17	-42.79	0.00	-42.38	-44.70	0.00	0.00	-39.77	-25.62	-22.80	-40.08	0.00	-30.02
3	-38.24	-32.12	-43.18	0.00	-42.91	-43.18	0.00	-48.01	-33.66	-28.79	-37.37	0.00	0.00	-32.47
4	-38.79	-32.82	-42.27	0.00	-41.20	-30.67	-40.94	-28.66	-28.00	-29.62	-38.21	0.00	0.00	-31.42
5	-34.93	-32.86	-46.04	0.00	-45.23	-37.68	-37.08	-35.45	-34.69	-28.45	-38.33	0.00	0.00	-31.05

Q-learning Policy Visualization

Below is the presentation of policy of the agent following the action sequence taken for the maximum reward, which is 642.

As observed shown below, the agent is travelling correctly with the action sequence = [4, 1, 1, 1, 3, 1, 4, 4, 4, 4, 1, 4, 2, 2, 2, 3, 2, 4, 4, 4, 4, 4, 2, 4, 1, 1, 1, 3, 1, 4, 4, 4, 1, 4, 2, 2, 2]

```
In [8]: env = CliffBoxGridWorld(render=True)
state = env.reset()
env.print_world()
done = False
rewards = []
step = 0

while not done:
    action = agent.max_action_seq[step]
    state, reward, done, info = env.step(action)
    rewards.append(reward)
    print(f'step: {env.timesteps}, state: {state}, action taken: {agent.action_seq[action]}'
          env.print_world()
    step += 1

print(f'rewards: {sum(rewards)}')
print(f'action history: {env.action_history}')
```

```
[[b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'_' b'_'
  b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'_' b'_'
  b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'x' b'_'
  b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'x' b'x' b'_'
  b'_' b'B' b'_' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'G']
[b'A' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' ]]
```

step: 1, state: (5, 1, 4, 1), action taken: 3, reward: -14

Action: 4

```
[[b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'_' b'_'
  b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'_' b'_'
  b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'x' b'_'
  b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'x' b'x' b'_'
  b'_' b'B' b'_' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'G']
[b'_' b'A' b'_' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' ]]
```

step: 2, state: (4, 1, 3, 1), action taken: 1, reward: -15

Action: 1

```
[[b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'_' b'_'
  b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'_' b'_'
  b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'x' b'_'
  b'_' b'B' b'_' b'x' b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'x' b'x' b'_'
  b'_' b'A' b'_' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'G']
[b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' ]]
```

step: 3, state: (3, 1, 2, 1), action taken: 1, reward: -16

Action: 1

```
[[b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'_' b'_'
  b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'_' b'_'
  b'_' b'B' b'_' b'x' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'x' b'_'
  b'_' b'A' b'_' b'x' b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'x' b'x' b'_'
  b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'G']
[b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' ]]
```

step: 4, state: (2, 1, 1, 1), action taken: 1, reward: -17

Action: 1

```
[[b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'_' b'_'
  b'_' b'B' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'_' b'_'
  b'_' b'A' b'_' b'x' b'_' b'_' b'_' b'x' b'x' b'_' b'_' b'_' b'_' b'x' b'_'
  b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'x' b'x' b'_'
  b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'G']
```


[illegible]

[[b' ' b' ' b' ' b' ' b' ' b' ' b'x' b'x' b' ' b' ' b' ' b' ' b' ']


```
[b'_' b'_' b'_' b'x' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'_' b'x' b'x' b'_']]  
rewards: 642  
action history: [4, 1, 1, 1, 3, 1, 4, 4, 4, 4, 1, 4, 2, 2, 2, 3, 2, 4, 4, 4, 4, 4, 2, 4,  
1, 1, 1, 3, 1, 4, 4, 4, 1, 4, 2, 2, 2]
```