## CarMDP Environment, which is provided for your convenience. You should not change code of this environment. This Jupyter notebook is prepared by Kui Wu

You have implemented a model-based solution in Assignment 2, and a model-free solution in Assignment 3. In this assignment, we assume that the agent initially knows nothing about the environment but gradually learns the environment and uses the learned knowledge for planning. We build the solution based on the Dyna framework.

```
In [1]:
        import numpy as np
        import matplotlib.pyplot as plt
        import gym
        import random
        from gym import Env
        class CarMDP(Env):
            11 11 11
            Car MDP with simple stochastic dynamics.
            The states are tuples with two elements:
                - a position index (i, j)
                - an integer from (0, 1, 2, 3) representing absolute orientation (se
            For example, the state
                s = (0, 1, 2)
            represents the car in the cell with indices (0, 1) and oriented to face
            def __init__(self, width, height, obstacles, goal transition, initial st
                         collision reward=-5., goal reward=10., stagnation penalty=
                self.width = width
                self.height = height
                self.grid map = np.ones((width, height))
                for cell in obstacles:
                     self.grid_map[cell[0], cell[1]] = 0.
                self.obstacles = obstacles
                self.orientations = {0: 'North', 1: 'East', 2: 'South', 3: 'West'}
                self.A = {0: 'Forward', 1: 'Left', 2: 'Right', 3: 'Brake'}
                self.goal transition = goal transition # Tuple containing start and
                self.p corr = p corr
                self.p_err = (1. - p_corr)/2.
                self.base reward = base reward
                self.collision reward = collision reward
                self.goal_reward = goal_reward
                self.stagnation penalty = stagnation penalty
                self.state_history = []
                self.action history = []
                self.reward history = []
                assert initial state[0] >= 0 and initial state[1] >= 0 and initial state
                        initial state[2] in self.orientations, "ERROR: initial state
                self.state_history = [initial_state]
                self.action history = []
                self.reward_history = []
                self.init state=initial state
            def reset(self):
                self.state history = [self.init state]
                self.action history = []
                self.reward history = []
            def is collision(self, state):
                is out of bounds = state[0] < 0 or state[0] >= self.width or state[]
                      state[1] >= self.height
                return is out of bounds or (state[0], state[1]) in self.obstacles
```

```
def transition dynamics(self, state, action):
   assert not self.is collision(state), "ERROR: can't take an action for
   delta = 1
   orientation = state[2]
    if self.orientations[orientation] == 'North':
        left = (state[0] - delta, state[1] - delta)
        forward = (state[0], state[1] - delta)
        right = (state[0] + delta, state[1] - delta)
    elif self.orientations[orientation] == 'West':
        left = (state[0] - delta, state[1] + delta)
        forward = (state[0] - delta, state[1])
        right = (state[0] - delta, state[1] - delta)
    elif self.orientations[orientation] == 'South':
        left = (state[0] + delta, state[1] + delta)
        forward = (state[0], state[1] + delta)
        right = (state[0] - delta, state[1] + delta)
    elif self.orientations[orientation] == 'East':
       left = (state[0] + delta, state[1] - delta)
        forward = (state[0] + delta, state[1])
        right = (state[0] + delta, state[1] + delta)
    # p gives categorical distribution over (state, left, forward, righ
    if self.A[action] == 'Forward':
       p = np.array([0., self.p_err, self.p_corr, self.p_err])
    elif self.A[action] == 'Right':
        p = np.array([0., 0., 2.*self.p_err, self.p_corr])
    elif self.A[action] == 'Left':
        p = np.array([0., self.p corr, 2. * self.p err, 0.])
    elif self.A[action] == 'Brake':
        p = np.array([self.p corr, 0., 2. * self.p err, 0.])
    candidate next state positions = (state, left, forward, right)
    next state position = candidate next state positions[categorical sar
    # Handle orientation dynamics (deterministic)
   new orientation = orientation
    if self.A[action] == 'Right':
       new orientation = (orientation + 1) % 4
    elif self.A[action] == 'Left':
       new orientation = (orientation - 1) % 4
    return next state position[0], next state position[1], new orientat
def step(self, action):
   assert action in self.A, f"ERROR: action {action} not permitted"
    terminal = False
    current state = self.state history[-1] # -1 means the current eleme
   next state = self.transition dynamics(current state, action)
    if self.is collision(next state):
        reward = self.collision reward
        terminal = True
    elif (current state[0], current state[1]) == self.goal transition[0
            (next state[0], next state[1]) == self.goal transition[1]:
        reward = self.goal_reward
        terminal = True  # TODO: allow multiple laps like this?
    elif current state == next state:
        reward = self.stagnation penalty
```

```
terminal = False
        else:
            reward = self.base reward
            terminal = False
        self.state history.append(next state)
        self.reward history.append(reward)
        self.action history.append(action)
        return next state, reward, terminal, []
    def render(self, title):
        self. plot history(title)
    def plot history(self, title):
        Plot the MDP's trajectory on the grid map.
        :param title:
        :return:
        11 11 11
        fig = plt.figure()
        plt.imshow(self.grid map.T, cmap='gray')
        plt.grid()
        x = np.zeros(len(self.state history))
        y = np.zeros(x.shape)
        for idx in range(len(x)):
            x[idx] = self.state history[idx][0]
            y[idx] = self.state history[idx][1]
            if self.state history[idx][2] == 0:
                plt.arrow(x[idx], y[idx], 0., -0.25, width=0.1)
            elif self.state history[idx][2] == 1:
                plt.arrow(x[idx], y[idx], 0.25, 0., width=0.1)
            elif self.state history[idx][2] == 2:
                plt.arrow(x[idx], y[idx], 0., 0.25, width=0.1)
            else:
                plt.arrow(x[idx], y[idx], -0.25, 0., width=0.1)
        plt.plot(x, y, 'b-') # Plot trajectory
        plt.xlim([-0.5, self.width + 0.5])
        plt.ylim([self.height + 0.5, -0.5])
        plt.title(title)
        plt.xlabel('x')
        plt.ylabel('y')
        plt.show()
        return fig
def categorical sample index(p: np.ndarray) -> int:
    Sample a categorical distribution.
    :param p: a categorical distribution's probability mass function (i.e.,
              returning idx for an integer 0 <= idx < len(p)). I.e., np.sum
    :return: index of a sample weighted by the categorical distribution des
    11 11 11
    P = np.cumsum(p)
    sample = np.random.rand()
    return np.argmax(P > sample)
```

Below is the skeleton code of your agent. Your solution should be filled here. You may need to introduce new functions and/or new data structure here.

```
In [2]:
        class ReinforcementLearningAgent:
            Your implementation of a reinforcement learning agent.
            Feel free to add additional methods and attributes.
            def init (self):
                ### STUDENT CODE GOES HERE
                # Set any parameters
                \# You can add arguments to init , so long as they have default v_0
                self.qamma = 1
                self.epsilon = 0.05
                self.alpha = 0.05
                self.Q = np.zeros((9*6, 4))
                self.shape = (9,6)
                self.model = {}
                self.state history = []
                self.action_history = []
            def reset(self, init_state) -> int:
                Called at the start of each episode.
                :param init state:
                :return: first action to take.
                ### STUDENT CODE GOES HERE
                action = 0
                self.state history = []
                self.action history = []
                self.state history.append(init state)
                state = (init state[0], init state[1])
                state = np.ravel multi index(tuple(state), self.shape)
                if np.random.random() < self.epsilon:</pre>
                     action = np.random.randint(4)
                else:
                     action = np.argmax(self.Q[state,:])
                self.action history.append(action)
                return action
            def next action(self, reward: float, state: int, terminal: bool) -> int
                Called during each time step of a reinforcement learning episode
                :param reward: reward resulting from the last time step's action
                :param state: state resulting from the last time step's action
                :param terminal: bool indicating whether state is a terminal state
                :return: next action to take
                ### STUDENT CODE GOES HERE
                # Produce the next action to take in an episode as a function of the
                # You may find it useful to track past actions, states, and rewards
                # Additionally, algorithms that learn during an episode (e.g., tempo
                  if terminal:
```

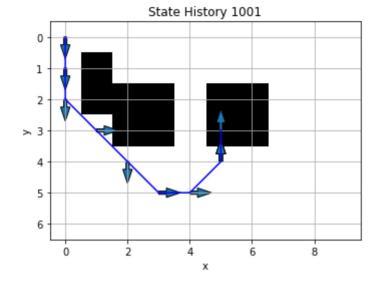
```
print(state)
#
             return 0
       prev state = self.state history[-1]
       prev state = (prev state[0], prev state[1])
       prev action = self.action history[-1]
       prev state = np.ravel multi index(tuple(prev state), self.shape)
        if terminal:
            next state = -1
            self.Q[prev state, prev action] += self.alpha*(reward - self.Q[]
            self.model learning(prev state, prev action, next state, reward
            self.planning(5)
            return 0
        self.state_history.append(state)
        next state = (state[0], state[1])
       next state = np.ravel multi index(tuple(next state), self.shape)
       eps action = 0
       if np.random.random() < self.epsilon:</pre>
           eps action = np.random.randint(4)
        else:
            eps_action = np.argmax(self.Q[next state,:])
        self.action history.append(eps action)
         prev q value = self.Q[prev state, prev action]
         next q value = self.Q[next state, eps action]
         prev q value += self.alpha*(reward + self.gamma*next q value - prediction)
        self.Q[prev state, prev action] += self.alpha*(reward + self.gamma*;
        self.model_learning(prev_state, prev_action, next_state, reward)
        self.planning(5)
        #print(self.model)
        #print(self.Q)
        return eps action
    def planning(self, planningStep):
        ### STUDENT CODE GOES HERE
        # Set any parameters
        # You can add other arguments to planning
        for i in range(planningStep):
            state action, 1 = random.choice(list(self.model.items()))
            prev state = state action[0]
            prev action = state action[1]
            total count = 0
            for j in 1:
                total count += j[2]
            # Do an expected update instead of sample update
            update = 0
            for k in 1:
                next state = k[0]
                reward = k[1]
                freq = k[2]
```

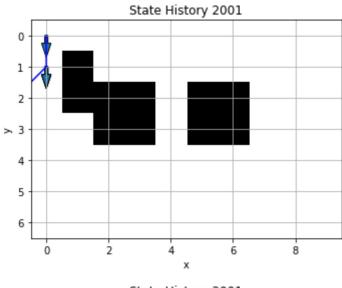
```
prob = freq/total count
            if next state == -1:
                value = prob*reward
                update += value
            else:
                next action = 0
                if np.random.random() < self.epsilon:</pre>
                    next action = np.random.randint(4)
                    next action = np.argmax(self.Q[next state,:])
                value = prob*(reward + self.gamma*self.Q[next state,nex
            #value = prob*(reward + self.gamma*np.max(self.Q[next state
                update += value
        self.Q[prev state,prev action] = update
def model learning(self, prev state, prev action, next state, reward):
    ### STUDENT CODE GOES HERE
    # Set any parameters
    # You can add other arguments to model learning
    # model stores the state-action pair as key and [next state, reward
    state action = tuple((prev state,prev action))
    if state_action not in self.model:
        l = [[next state, reward, 1]]
        self.model.update({state action : 1})
        l = self.model[state action]
        found = 0
        for i in 1:
            if i[0] == next state:
                found = 1
                i[1] = i[1] + (1/i[2]) * (reward-i[1])
                i[2] += 1
                break
        if found == 0:
            temp = [next state,reward,1]
            1.append(temp)
            #print("hi")
        self.model.update({state action : 1})
def finish_episode(self):
   Called at the end of each episode.
    :return: nothing
    ### STUDENT CODE GOES HERE
    # Algorithms that learn from an entire episode (e.g., Monte Carlo)
   pass
```

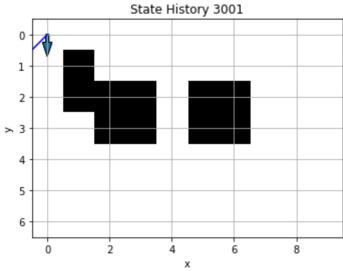
Below is the sample test code. In the final print out you need to print out the correct policy name (It is random so far). Note that

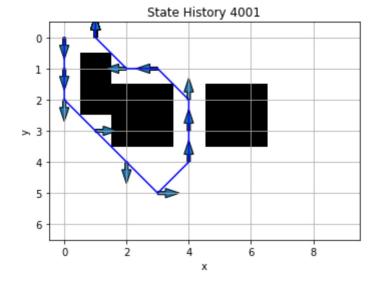
```
DOILV HATHE TILLS LAHAOHE SO LALL. LYOTE THAT
In [3]:
        def test rl algorithm(rl agent, car mdp, initial state, n episodes=10000, n
            Code that will be used to test your implementation of ReinforcementLear
            As you can see, you are responsible for implementing three methods in Re
                - reset (called at the start of every episode)
                - next action (called at every time step of an episode)
                - finish episode (called at the end of each episode)
            :param rl_agent: an instance of your ReinforcementLearningAgent class
            :param car mdp: an instance of CarMDP
            :param init state: the initial state
            :param n episodes: number of episodes to use for this test
            :param n_plot: display a plot every n_plot episodes
            :return:
            11 11 11
            returns = []
            for episode in range(n episodes):
                G = 0. # Keep track of the returns for this episode (discount factor)
                # Re-initialize the MDP and the RL agent
                car mdp.reset();
                action = rl agent.reset(initial state)
                terminal = False
                while not terminal: # Loop until a terminal state is reached
                    next state, reward, terminal, [] = car mdp.step(action)
                    G += reward
                    action = rl agent.next action(reward, next state, terminal)
                rl agent.finish episode()
                returns += [G]
                # Plot the trajectory every n plot episodes
                if episode % n plot == 0 and episode > 0:
                    #print(n plot)
                    #print(episode)
                    car mdp.render('State History ' + str(episode + 1))
            return returns
        if name == '__main__':
            # Size of the CarMDP map (any cell outside of this rectangle is a termi
            width = 9
            height = 6
            initial state = (0, 0, 2) # Top left corner (0, 0), facing "Down" (2)
            obstacles = [(2, 2), (2, 3), (3, 2), (3, 3), # Cells filled with obstacles
                         (5, 2), (5, 3), (6, 2), (6, 3),
                         (1, 1), (1, 2)]
            goal\_transition = ((1, 0), (0, 0)) # Transitioning from cell (1, 0) to
            p corr = 0.95 # Probability of actions working as intended
          # Create environment
            car mdp = CarMDP(width, height, obstacles, goal transition, initial state
         # Create RL agent. # You must complete this class in your solution, it is
         # the first agent (rl agent) is just to track the agent to see if it is le
            rl agent = ReinforcementLearningAgent()
```

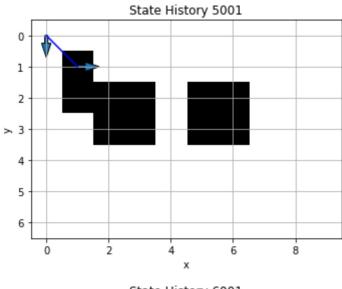
```
student returns = test rl algorithm(rl agent, car mdp, initial state, n
# Example plot. You need to change it according to the assignment requi
n runs = 10
n = 10000
returns = np.zeros((n runs, n episodes))
for run in range(n runs):
    rl agentnew = ReinforcementLearningAgent()
    returns[run, :] = test rl algorithm(rl agentnew, car mdp, initial s
# Plot one curve like this for each parameter setting - the template co
# returns a random action, so this example curve will just be noise. Wh
# should increase as the number of episodes increases. Feel free to chall
rolling average width = 100
# Compute the mean (over n runs) for each episode
mean return = np.mean(returns, axis=0)
# Compute the rolling average (over episodes) to smooth out the curve
rolling average mean return = np.convolve(mean return, np.ones(rolling a
plt.figure()
plt.plot(rolling average mean return, 'b-') # Plot the smoothed average
plt.grid()
plt.title('Learning Curve')
plt.xlabel('Episode')
plt.ylabel('Average Return')
plt.legend(['alpha = 0.05, epsilon = 0.05, planning steps = 5'])
plt.show()
```

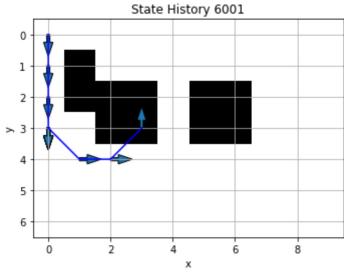


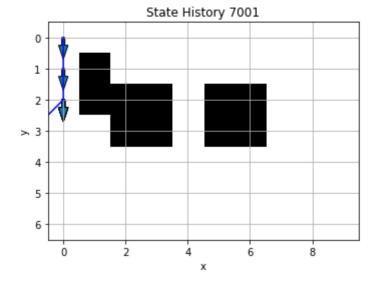


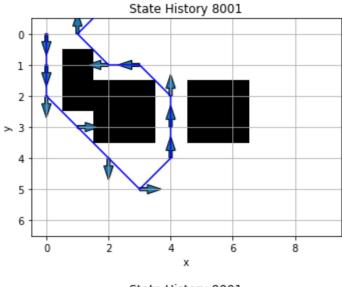


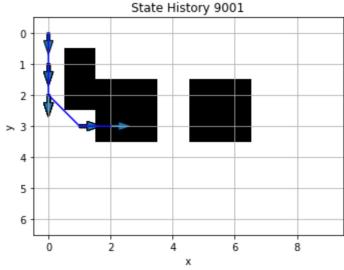


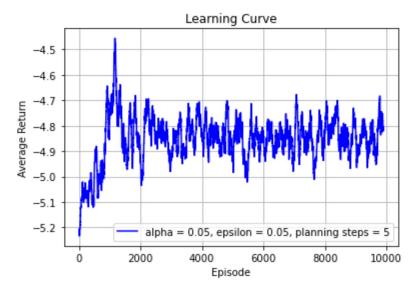












In [ ]:

In [ ]: