**Giorgio Mendoza**

**RBE595-S24-S04**

**DP Programming Excercise**

Description: The code aims to illustrate how different reinforcement learning algorithms (Policy Iteration, Generalized Policy Iteration, Value Iteration) can be applied to find optimal navigation strategies in a grid world. It also demonstrates both deterministic and stochastic behaviors, showcasing how the agent's decision-making changes under uncertainty.

The agent can start from any free space within the grid world. Free spaces are represented by cells in the WorldGrid that have a value of 0.

The goal location is specified by the desired_position variable, which is initialized in the main execution part of the code. In this specific implementation, the goal is set at the coordinates (7, 10)

Objective: The agent's objective is to find the most efficient path to the goal while avoiding obstacles. This is achieved by learning an optimal policy through different reinforcement learning algorithms.

The learned policy guides the agent on which action to take in each state (position in the grid) to maximize its cumulative reward, which involves successfully reaching the goal and minimizing collisions with obstacles.

```python
1  import numpy as np
2  from itertools import product
3  import matplotlib.pyplot as plt
4  from google.colab import drive
5  drive.mount('/content/drive')
6
7  reward_obstacle_collision_init = -50
8  reward_desired_goal_init = 100
9  tolerance_init = 1
10 mode_init = True
11 actions_init = 8
12 discount_factor_init = 0.95
13 reward_clear_move_init = -1.0
14 desired_position_init = (7,10)
15
16 class Agent:
17     def __init__(self, WorldGrid, total_actions = actions_init, discount_factor = discount_factor_init, desired_position = desired_po
18                  reward_obstacle_collision = reward_obstacle_collision_init, reward_desired_goal = reward_desired_goal_init, toler
19         self.total_actions = total_actions  # total number of actions
20         self.discount_factor = discount_factor  # discount factor for future rewards
21         self.desired_position = desired_position  # location of the goal
22         self.reward_clear_move = reward_clear_move  # reward for a clear move
23         self.reward_desired_goal = reward_desired_goal  # reward for reaching the goal
24         self.reward_obstacle_collision = reward_obstacle_collision  # penalty for hitting an obstacle
25         self.tolerance = tolerance  # convergence threshold
26         self.deterministic = mode  # flag for deterministic or stochastic model
27         self.WorldGrid = WorldGrid  # world definition
28         self.widthGrid, self.heightGrid = self.WorldGrid.shape  # dimensions of the world
29         # filter out boundary and obstacle states
30         self.states = [state for state in product(range(self.widthGrid), range(self.heightGrid)) if self.WorldGrid[state] == 0]
31         self.valid_states = len(self.states)  # total number of valid states
32         print("Total number of valid states:", self.valid_states)
33
34         # define movement actions and their state transitions
35         self.valid_movement_states = {
36             "up": (-1, 0),
37             "down": (1, 0),
38             "right": (0, 1),
39             "left": (0, -1),
40             "up_right": (-1, 1),
41             "up_left": (-1, -1),
42             "down_right": (1, 1),
43             "down_left": (1, -1)
44         }
45         # define actions with their probabilities
46         if self.deterministic:
47             self.valid_actions = {key: [(key, 1)] for key in self.valid_movement_states}
48         else:
49             # stochastic model with 60% intended move, 20% each for adjacent moves
50             self.valid_actions = {
51                 "up": [("up", 0.6), ("up_left", 0.2), ("up_right", 0.2)],
52                 "up_left": [("up_left", 0.6), ("up", 0.2), ("left", 0.2)],
53                 "up_right": [("up_right", 0.6), ("up", 0.2), ("right", 0.2)],
54                 "down": [("down", 0.6), ("down_left", 0.2), ("down_right", 0.2)],
55                 "down_left": [("down_left", 0.6), ("down", 0.2), ("left", 0.2)],
56                 "down_right": [("down_right", 0.6), ("down", 0.2), ("right", 0.2)],
57                 "left": [("left", 0.6), ("up_left", 0.2), ("down_left", 0.2)],
58                 "right": [("right", 0.6), ("up_right", 0.2), ("down_right", 0.2)]
59             }
60
61         # function for summing state and action tuples
62         self.state_plus_action = lambda state, action: tuple(map(sum, zip(state, self.valid_movement_states[action])))
63
64         # initialize policy with equal probability for each action
65         self.init_policy = {state: {action: 1 / self.total_actions for action in self.valid_actions} for state in self.states}
66         self.init_value_function = np.zeros_like(self.WorldGrid)
67
68     def PolicyGraph(self, policy, string):
69         fig, ax = plt.subplots(figsize=(15, 8))  # Create figure and axes
70         ax.set_title(string)
71         goal_y_coordinate, goal_x_coordinate = self.desired_position
72         plt.plot(goal_x_coordinate + 0.5, goal_y_coordinate + 0.5, "ro", markersize=10)  # mark the goal with a red dot
73
74         # Draw grid lines
75         for i in range(self.WorldGrid.shape[0] + 1):
76             ax.axhline(i, lw=1, color='black', zorder=5)
77         for i in range(self.WorldGrid.shape[1] + 1):
78             ax.axvline(i, lw=1, color='black', zorder=5)
79
80         # Draw walls (obstacles)
81         for y in range(self.WorldGrid.shape[0]):
82             for x in range(self.WorldGrid.shape[1]):
```

```
 83                    if self.WorldGrid[y, x] == 1:
 84                        ax.add_patch(plt.Rectangle((x, self.WorldGrid.shape[0] - y - 1), 1, 1, fill=True, color='black', zorder=5))
 85
 86        # Add policy arrows
 87        for state in policy:
 88            if 1 in policy[state].values():
 89                y, x = state
 90                action = max(policy[state], key=policy[state].get)
 91                dy, dx = self.valid_movement_states[action]
 92                ax.arrow(x + 0.5, self.WorldGrid.shape[0] - y - 0.5, dx * 0.3, -dy * 0.3, head_width=0.2, head_length=0.2, fc='blue',
 93
 94
 95        # Set axis limits and remove labels
 96        ax.set_xlim(0, self.WorldGrid.shape[1])
 97        ax.set_ylim(0, self.WorldGrid.shape[0])
 98        ax.set_xticks([])
 99        ax.set_yticks([])
100        ax.invert_yaxis()  # Invert y-axis to match the matrix representation
101
102        plt.show()
103
104
105    def ValueGraph(self, value, string):
106        plt.figure(figsize=(8, 8))  # larger figure size for clarity
107        plt.title(string)
108        goal_y_coordinate, goal_x_coordinate = self.desired_position
109        plt.plot(goal_x_coordinate, goal_y_coordinate, "ro", markersize=10)  # mark the goal with a red dot
110        im = plt.imshow(value, cmap="hot", interpolation='none')  # 'hot' colormap for value representation
111
112        # add a colorbar to indicate value scale
113        plt.colorbar(im)
114        plt.show()
115
116    def PolicyEval(self, policy, value_function):
117        max_change = 0
118        for state in self.states:
119            # initialize the value for this state
120            val = 0
121
122            # loop over all actions and their probabilities
123            for action, action_probs in self.valid_actions.items():
124                pi = policy[state][action]
125
126                # sum the value for all possible outcomes of this action
127                for _a, prob in action_probs:
128                    next_state = self.state_plus_action(state, _a)
129
130                    # if next state is valid (not out of bounds or an obstacle)
131                    if next_state in self.states:
132                        SetReward = self.SetReward(next_state)
133                        val += pi * prob * (SetReward + self.discount_factor * value_function[next_state])
134
135            # update the maximum max_change and value function for this state
136            max_change = max(max_change, abs(val - value_function[state]))
137            value_function[state] = val
138
139        return value_function, max_change
140
141    def PolicyImprov(self, policy, value):
142        is_converged = True
143        # iterate over all states to improve policy
144        for state in self.states:
145            # find the best action according to the current value function
146            best_action_value = float('-inf')
147            best_action = None
148
149            # examine the value of each action
150            for action in self.valid_actions:
151                action_value = 0
152                # consider the outcome of each action
153                for _a, prob in self.valid_actions[action]:
154                    next_state = self.state_plus_action(state, _a)
155                    # calculate value if next state is valid
156                    if next_state in self.states:
157                        SetReward = self.SetReward(next_state)
158                        action_value += prob * (SetReward + self.discount_factor * value[next_state])
159
160                # update the best action if this action is better
161                if action_value > best_action_value:
162                    best_action_value = action_value
163                    best_action = action
164
```
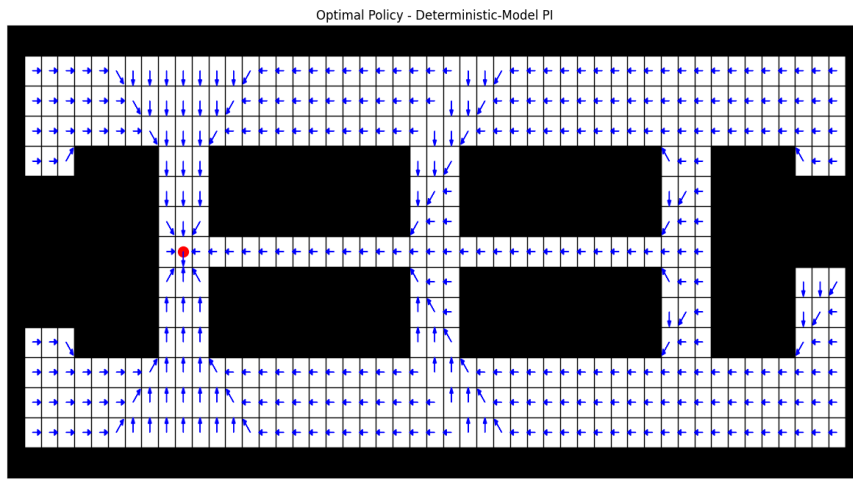
```
165              # compare the best action with the current policy's action
166              current_policy_action = max(policy[state], key=policy[state].get, default=None)
167              if current_policy_action != best_action:
168                  is_converged = False
169
170              # update the policy for this state
171              for action in self.valid_actions:
172                  policy[state][action] = 1 if action == best_action else 0
173
174          return policy, value, is_converged
175
176      def PolicyIter(self):
177          # initialize policy and value function
178          policy = dict(self.init_policy)
179          value = np.array(self.init_value_function)
180
181          # loop for policy iteration
182          while True:
183              # policy evaluation - only needs to run until max_change < tolerance once per iteration
184              value, max_change = self.PolicyEval(policy, value)
185              while max_change >= self.tolerance:
186                  value, max_change = self.PolicyEval(policy, value)
187
188              # policy improvement step
189              policy, value, is_converged = self.PolicyImprov(policy, value)
190
191              # if policy is stable, we're done
192              if is_converged:
193                  break
194
195          return policy, value
196
197      def GPI(self):
198          # initialize policy and value function
199          policy = dict(self.init_policy)
200          value = np.array(self.init_value_function)
201
202          # run generalized policy iteration
203          while True:
204              # policy evaluation with a single pass
205              value, _ = self.PolicyEval(policy, value)
206
207              # policy improvement step
208              policy, value, is_converged = self.PolicyImprov(policy, value)
209
210              # if policy is stable, the iteration stops
211              if is_converged:
212                  break
213
214          return policy, value
215
216      def ValueIter(self, plot=False):
217          value = np.array(self.init_value_function)
218
219          # keep iterating until no significant changes are made to the value function
220          while True:
221              max_change = 0
222              for state in self.states:
223                  # compute the value of each action and select the best one
224                  best_action_value = float('-inf')
225                  for action, action_probs in self.valid_actions.items():
226                      action_value = sum(
227                          prob * (self.SetReward(self.state_plus_action(state, _a)) + self.discount_factor * value[self.state_plus_act
228                          for _a, prob in action_probs if self.state_plus_action(state, _a) in self.states
229                      )
230                      best_action_value = max(best_action_value, action_value)
231
232                  # track the largest change from the current value function
233                  max_change = max(max_change, abs(best_action_value - value[state]))
234                  # update the value function with the best action value
235                  value[state] = best_action_value
236
237              # display the current value function if live plotting is enabled
238              if plot:
239                  self.ValueGraph(value, "Value Iteration Live Value Plot")
240
241              # stop if the change is below the threshold for all states
242              if max_change < self.tolerance:
243                  break
244
245          # construct a policy where each state takes the action leading to the highest value
246          policy = {state: {action: 0 for action in self.valid_actions} for state in self.states}
```
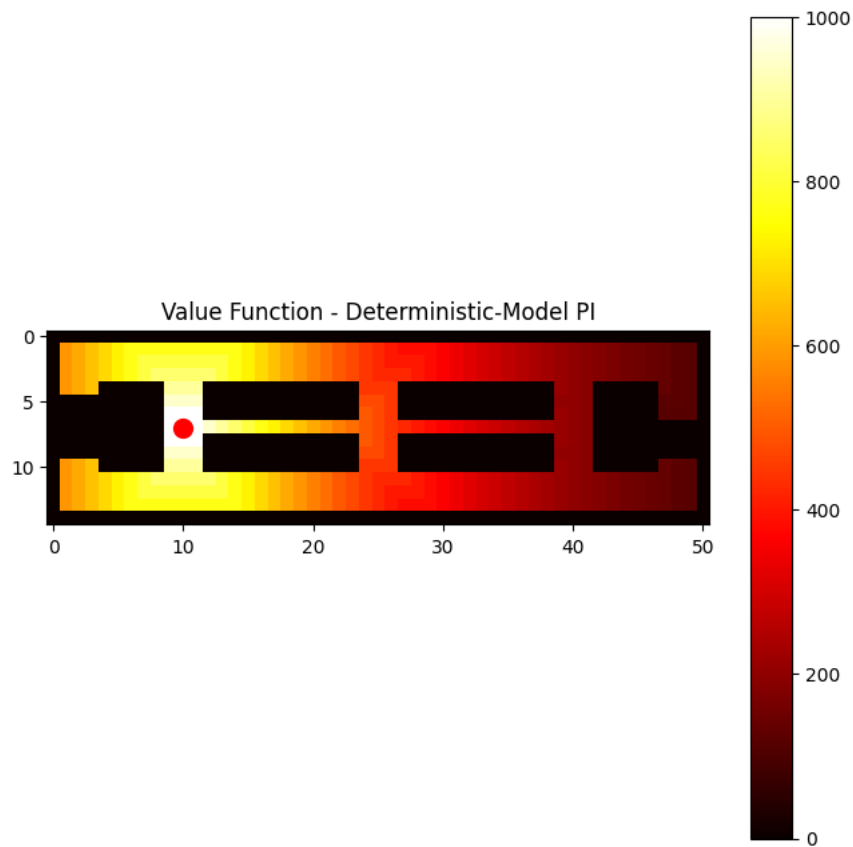
```
247        for state in self.states:
248            _, best_action = max(
249                ((sum(prob * (self.SetReward(self.state_plus_action(state, _a)) + self.discount_factor * value[self.state_plus_action
250                   for _a, prob in action_probs if self.state_plus_action(state, _a) in self.states), action)
251                 for action, action_probs in self.valid_actions.items()),
252                key=lambda x: x[0]
253            )
254            policy[state][best_action] = 1
255
256        return policy, value
257
258    def SetReward(self, state):
259        # returns the reward for the robot's state after an action
260        # hitting an obstacle yields reward_obstacle_collision
261        # reaching the goal yields reward_desired_goal
262        # any other move yields reward_clear_move
263
264        if state == self.desired_position:
265            return self.reward_desired_goal
266        elif self.WorldGrid[state] == 1:  # Direct indexing since state is a tuple
267            return self.reward_obstacle_collision
268        return self.reward_clear_move
269
270 if __name__ == "__main__":
271     # define the world as a numpy array
272     WorldGrid = np.array(
273 [
274 [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
275 [1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
276 [1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
277 [1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
278 [1, 0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1,
279 [1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1,
280 [1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1,
281 [1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1,
282 [1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1,
283 [1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1,
284 [1, 0, 0, 0, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 1,
285 [1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
286 [1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
287 [1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
288 [1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
289
290  # adjust goal coordinates indexing
291      # create agent instance
292      agent = Agent(WorldGrid, desired_position=(7, 10))

      Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).
      Total number of valid states: 399
```

The collective arrows form paths leading towards a common destination. The arrows seem to avoid black squares, adhering to the constraints of the environment. This suggests that the policy has successfully learned to navigate around obstacles. As shown below, the arrows converge towards the goal state in the direction where most arrows in the free spaces point towards. Since this is a "Deterministic-Model PI", each state has one clear action to take (one arrow per state), which aligns with a deterministic approach where the outcome of an action is certain.

```
1     # perform policy iteration for deterministic model
2     policy, value = agent.PolicyIter()
3     agent.PolicyGraph(policy, "Optimal Policy - Deterministic-Model PI")
```

Optimal Policy - Deterministic-Model PI

This plot visualizes the value function for the deterministic model using policy iteration. The color gradient, ranging from yellow to dark red, represents the value at each state in the grid, with higher values in yellow and lower values in dark red. The filled red circle indicates the goal position, which is the state with the highest value—this is consistent with the idea that reaching the goal yields the highest reward.

Black squares represent obstacles or walls where the agent cannot go. The intensity of the colors reflects the potential value of being in a particular state; states closer to the goal tend to have higher values since they are closer to achieving the reward. This value function helps to guide the agent's decisions: at each state, the agent will choose the action that leads to the state with the highest value.
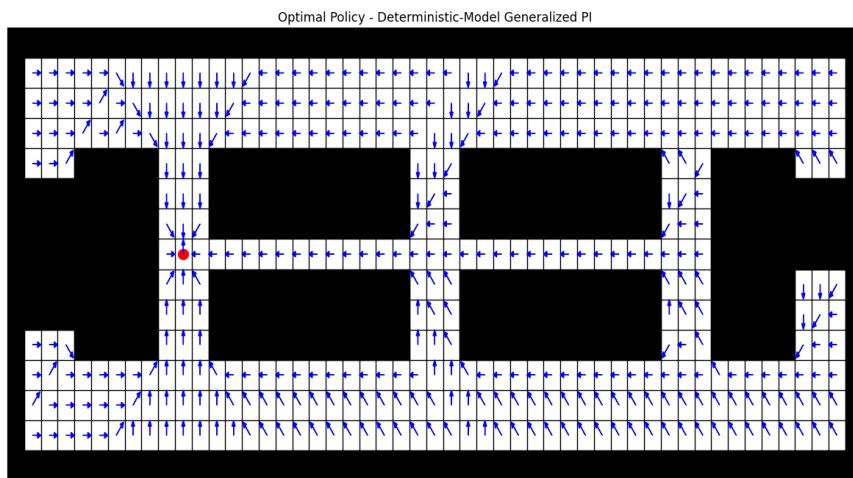
```
1    agent.ValueGraph(value, "Value Function - Deterministic-Model PI")
```

Value Function - Deterministic-Model PI

```
1    # perform generalized policy iteration for deterministic model
2    policy, value = agent.GPI()
3    agent.PolicyGraph(policy, "Optimal Policy - Deterministic-Model Generalized PI")
```
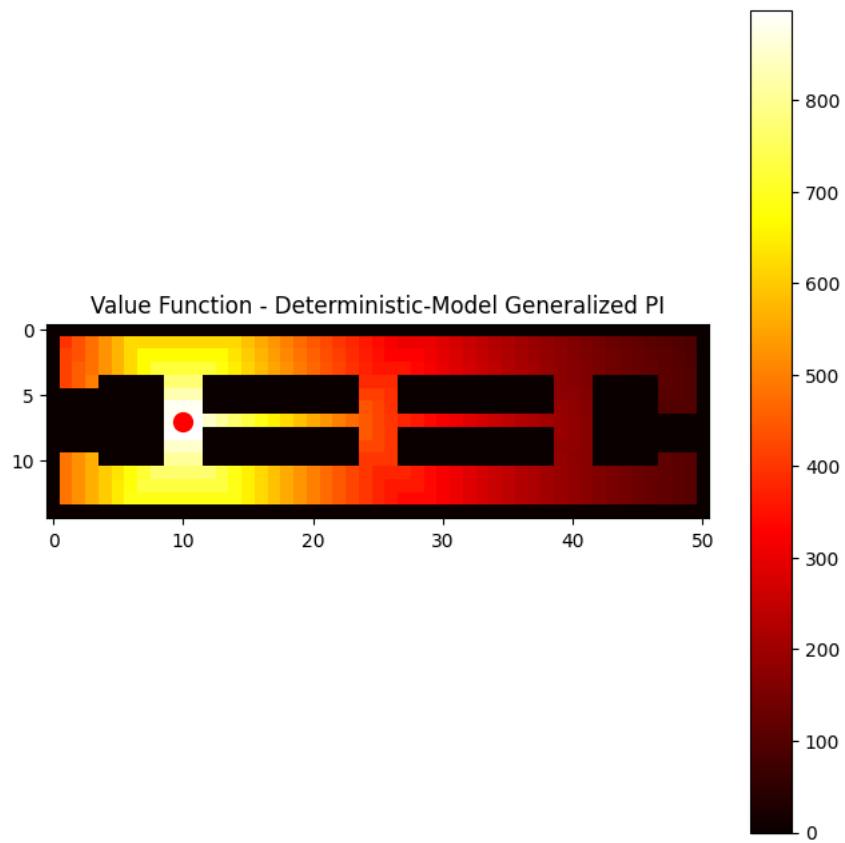


Optimal Policy - Deterministic-Model Generalized PI

```
1    agent.ValueGraph(value, "Value Function - Deterministic-Model Generalized PI")
```
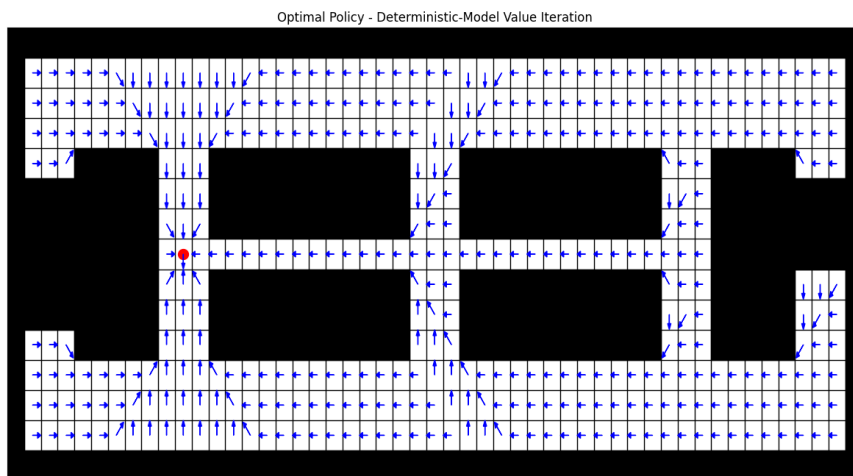
Value Function - Deterministic-Model Generalized PI

```
1    # perform value iteration for deterministic model
2    policy, value = agent.ValueIter()
3    agent.PolicyGraph(policy, "Optimal Policy - Deterministic-Model Value Iteration")
```



Optimal Policy - Deterministic-Model Value Iteration

```
1    agent.ValueGraph(value, "Value Function - Deterministic-Model Value Iteration")
```

Value Function - Deterministic-Model Value Iteration