

Stochastic Wind GridWorld King's Moves

Exercise 6.10: King's Move

Based from the plot that I was able to produce, it can be seen that with the fixed epsilon, the performance using SARSA is different from the other two methods. However, it can be seen that once the Variable Epsilon is used, the performances are roughly the same for all the different methods.

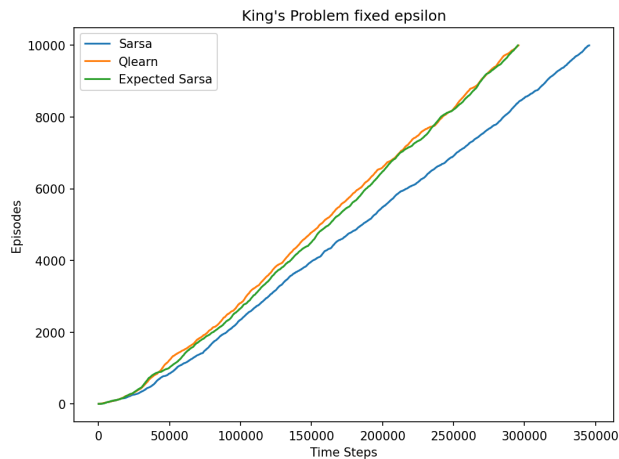


Figure 1: Fixed Epsilon

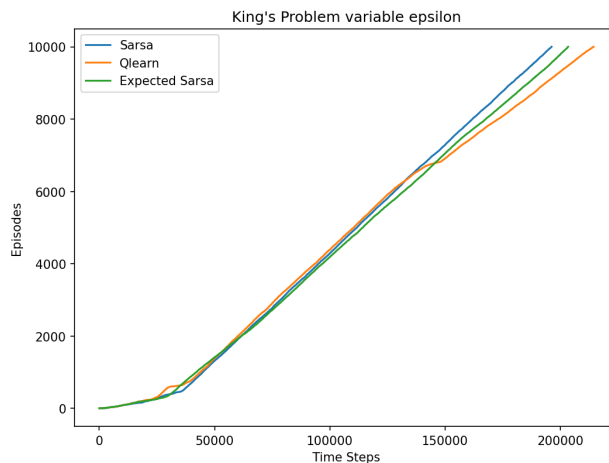


Figure 2: Variable Epsilon

Appendix

Python

```
1 '''
2 Justine Serdoncillo
3 IE 5571 - Dynamic Programming
4 HW 4 Exercise 6.10
5 November 6, 2023
6 '''
7
8 """
9 Exercise 6.10: Stochastic Wind (programming) Re-solve the windy gridworld task with
10 Kings moves, assuming that the effect of the wind, if there is any, is stochastic,
    sometimes
11 varying by 1 from the mean values given for each column. That is, a third of the time
12 you move exactly according to these values, as in the previous exercise, but also a third
13 of the time you move one cell above that, and another third of the time you move one
14 cell below that. For example, if you are one cell to the right of the goal and you move
15 left, then one-third of the time you move one cell above the goal, one-third of the time
16 you move two cells above the goal, and one-third of the time you move to the goal.
17 """
18
19 import numpy as np
20 import matplotlib.pyplot as plt
21 import random
22 import math
23
24 # %%
25 wind = [0, 0, 0, 1, 1, 1, 2, 2, 1, 0]
26
27 class Grid:
28     def __init__(self, w, h, wind):
29         self.w = w
30         self.h = h
31         self.wind = wind
32         self.actions = list(range(9))
33         self.states = [tuple([i, j]) for j in range(self.h) for i in range(self.w)]
34         self.starting_state = tuple([0, 3])
35         self.terminal_states = [tuple([7, 3])]
36
37     def take_action(self, state, action):
38         #locations are x, y
39         x = state[0]
40         y = state[1]
41         win_rand = random.random()
42         if win_rand <= 1/3:
43             y += wind[x]
44         elif win_rand <= 2/3:
45             y += wind[x] + 1
46         else:
47             y += wind[x] - 1
48         if action == 0: # up
49             x += 0
50             y += 1
51         if action == 1: # up, right
52             x += 1
53             y += 1
54         if action == 2: # right
55             x += 1
56             y += 0
57         if action == 3: # right, down
58             x += 1
59             y += -1
60         if action == 4: # down
61             x += 0
62             y += -1
63         if action == 5: # left, down
64             x += -1
65             y += -1
66         if action == 6: # left
```

```

67         x += -1
68         y += 0
69         if action == 7: # up, left
70             x += -1
71             y += 1
72         if action == 8: # No move
73             x += 0
74             y += 0
75         if x >= self.w:
76             x = self.w - 1
77         if y >= self.h:
78             y = self.h - 1
79         if y < 0:
80             y = 0
81         if x < 0:
82             x = 0
83         r = -1
84         return tuple([x, y]), r
85
86 # %% SARSA
87 class TDAIlg:
88     def __init__(self, problem, qs=None, policy=None, eps=.05, gamma=1, alpha=1, max_time=
None, special=False):
89         self.problem = problem
90         if qs is None:
91             self.qs = {s: {a: 0 for a in problem.actions} for s in problem.states}
92         else:
93             self.qs = qs
94         if policy is None:
95             self.policy = {s: random.choice(problem.actions) for s in problem.states}
96         else:
97             self.policy = policy
98         self.eps = eps
99         self.gamma = gamma
100        self.alpha = alpha
101        self.max_time = max_time
102        self.state = None
103        self.ep = 1
104        self.time = 0
105        self.ep_log = []
106        self.r_log = []
107        self.special = special
108
109    def get_action(self, state):
110        if self.special is not False:
111            self.eps = self.special/self.ep
112            if random.random() < self.eps:
113                choice = random.choice(self.problem.actions)
114            else:
115                choice = self.best_action(state)
116        return choice
117
118    def best_action(self, state):
119        choices = [a for a in self.problem.actions if self.qs[state][a] == max(self.qs[state]
).values()]]
120        best_action = random.choice(choices)
121        return best_action
122
123    def run_episode(self):
124        if (self.max_time is not None and self.time <= self.max_time) or self.max_time is
None:
125            old_count = self.time
126            cum_reward = self.update()
127            self.ep_log.extend([self.ep for _ in range(self.time - old_count)])
128            self.ep += 1
129            if len(self.r_log) > 0:
130                self.r_log.append(self.r_log[-1] + cum_reward)
131            else:
132                self.r_log.append(cum_reward)
133        else:
134            pass
135

```

```

136     def run(self, episodes):
137         i = 0
138         while i <= episodes:
139             self.run_episode()
140             i += 1
141
142     class Sarsa(TDAlg):
143         def update(self):
144             cum_reward = 0
145             state = self.problem.starting_state
146             action = self.get_action(state)
147             while state not in self.problem.terminal_states:
148                 self.time += 1
149                 new_state, reward = self.problem.take_action(state, action)
150                 new_action = self.get_action(new_state)
151                 target = reward + self.gamma*self.qs[new_state][new_action] - self.qs[state][
action]
152                 cum_reward += reward
153                 self.qs[state][action] += self.alpha*target
154                 action = new_action
155                 state = new_state
156             return cum_reward
157
158     # %% Q-Learning
159
160     class QLearn(TDAlg):
161         def update(self):
162             cum_reward = 0
163             state = self.problem.starting_state
164             while state not in self.problem.terminal_states:
165                 self.time += 1
166                 action = self.get_action(state)
167                 new_state, reward = self.problem.take_action(state, action)
168                 cum_reward += reward
169                 best = self.best_action(new_state)
170                 target = reward + self.gamma*self.qs[new_state][best] - self.qs[state][action]
171                 self.qs[state][action] += self.alpha*target
172                 state = new_state
173             return cum_reward
174
175         def update_from_model(self, s, a, s_, r):
176             best = self.best_action(s_)
177             target = r + self.gamma*self.qs[s_][best] - self.qs[s][a]
178             self.qs[s][a] += self.alpha*target
179
180     # %% Ex-pected SARSA
181
182     class ESarsa(TDAlg):
183         def update(self):
184             cum_reward = 0
185             state = self.problem.starting_state
186             while state not in self.problem.terminal_states:
187                 self.time += 1
188                 action = self.get_action(state)
189                 new_state, reward = self.problem.take_action(state, action)
190                 val = reward - self.qs[state][action]
191                 best = self.best_action(new_state)
192                 for a in self.problem.actions:
193                     if a == best:
194                         val += (1-self.eps)*self.gamma*self.qs[new_state][a]
195                         val += (self.eps/len(self.problem.actions))*self.gamma*self.qs[new_state][a]
196                 target = val
197                 cum_reward += reward
198                 self.qs[state][action] += self.alpha*target
199                 state = new_state
200             return cum_reward
201
202     # %% Fixed Epsilon
203     fig, ax = plt.subplots(figsize=(8,6), dpi=150)
204     ax.set_title("King's Problem fixed epsilon")
205     ax.set_xlabel("Time Steps")
206     ax.set_ylabel("Episodes")

```

```

207
208 grid = Grid(w=len(wind), h=7, wind=wind)
209 grid.actions = list(range(8))
210
211 sarsaFixed = Sarsa(grid, alpha=0.5, eps=.1)
212 sarsaFixed.run(10000)
213 ax.plot(sarsaFixed.ep_log, label='Sarsa')
214
215 QFixed = QLearn(grid, alpha=0.5, eps=.1)
216 QFixed.run(10000)
217 ax.plot(QFixed.ep_log, label='Qlearn')
218
219 EFixed = ESarsa(grid, alpha=0.5, eps=.1)
220 EFixed.run(10000)
221 ax.plot(EFixed.ep_log, label='Expected Sarsa')
222 ax.legend()
223
224 # %% Variable Epsilon
225 fig, ax = plt.subplots(figsize=(8,6), dpi=150)
226 ax.set_title("King's Problem variable epsilon")
227 ax.set_xlabel("Time Steps")
228 ax.set_ylabel("Episodes")
229
230 grid = Grid(w=len(wind), h=7, wind=wind)
231 grid.actions = list(range(8))
232
233 sarsaFixed = Sarsa(grid, alpha=0.5, eps=.3, special=.3)
234 sarsaFixed.run(10000)
235 ax.plot(sarsaFixed.ep_log, label='Sarsa')
236
237 QFixed = QLearn(grid, alpha=0.5, eps=.3, special=.3)
238 QFixed.run(10000)
239 ax.plot(QFixed.ep_log, label='Qlearn')
240
241 EFixed = ESarsa(grid, alpha=0.5, eps=.3, special=.3)
242 EFixed.run(10000)
243 ax.plot(EFixed.ep_log, label='Expected Sarsa')
244 ax.legend()

```

EXERCISE 5.6 WHAT IS THE EQUATION ANALOGOUS TO

$$V(s) = \frac{\sum_{t \in \tau(s)} R_t: T(t) - 1}{\sum_{t \in \tau(s)} R_t: T(t) - 1} G_t \quad \text{for action values } Q(s, a)$$

instead of state values $V(s)$, again given returns generated using b ?

now track state-action pairs instead then

$$Q(s, a) = \frac{\sum_{t \in \tau(s, a)} R_t: T(t) - 1}{\sum_{t \in \tau(s, a)} R_t: T(t) - 1} G_t$$

EXERCISE 5.8 THE RESULTS WITH EX 5.5 & SHOWN IN FIG 5.4 USED A FIRST VISIT MC.
CHANGE TO EVERY-VISIT MC, would variance of estimator $= \infty$? why/why not?

Estimator will still be infinite because reward is at terminal still.
& the visits to the state will still run to infinity, expected reward
at every state $= 1$

EXERCISE 6.7 DESIGN AN OFF-POLICY VERSION OF THE TD(0) UPDATE THAT CAN BE USED WITH ARBITRARY TARGET POLICY π & COVERING BEHAVIOURAL POLICY b , USING AT EACH STEP THE IMPORTANCE SAMPLING RATIO $P_{t:t} \quad (6.1)$

$$P_{t:T-1} = \frac{\prod_{k=t}^{T-1} \pi(A_k|S_k) p(S_{k+1}|S_k, A_k)}{\prod_{k=t}^{T-1} b(A_k|S_k) p(S_{k+1}|S_k, A_k)} = \prod_{k=t}^{T-1} \frac{\pi(A_k|S_k)}{b(A_k|S_k)}$$

$$P_{t:t} = \frac{\pi(A_t|S_t)}{b(A_t|S_t)}$$

TD update can be chosen as

(6.2)

$$V(S_t) = V(S_t) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(S_t)]$$

Can now make episode is made by b so update becomes

$$V_{\pi}(S_{t+1}) = V_{\pi}(S_{t+1}) + \alpha [P_{t:t} R_{t+1} + P_{t:t} V(S_{t+1}) - V(S_t)]$$

based on b

EXERCISE 6.2 SHOW THAT AN ACTION-VALUE VERSION OF (6.6)

$$G_t - V(S_t) = \sum_{k=t}^{T-1} \gamma^{k-t} \delta_k \quad \text{holds for the action-value form of the}$$

TD error $\delta_t = R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)$, again assuming that the

values don't change from step to step.

$$G_t - Q(S_t, A_t) = R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t) + \gamma Q(S_{t+1}, A_{t+1}) - \gamma Q(S_{t+1}, A_{t+1})$$

$$= \delta_t + \gamma [G_{t+1} - Q(S_{t+1}, A_{t+1})]$$

$$= \delta_t + \gamma [\delta_{t+1} + \gamma [G_{t+2} - Q(S_{t+2}, A_{t+2})]]$$

$$= \delta_t + \gamma \delta_{t+1} + \gamma^2 \{ \dots \}$$

$$= \sum_{k=t}^{T-1} \gamma^{k-t} \delta_k$$

$$G_t - Q(S_t, A_t) = \sum_{k=t}^{T-1} \gamma^{k-t} \delta_k$$

$$\delta_t = R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)$$

EXERCISE 6-12 SUPPOSE ACTION SELECTION IS GREEDY. IS Q-LEARNING THEN

EXACTLY THE SAME AS SARSA? WILL THEY MAKE EXACTLY THE SAME ACTION SELECTION & WEIGHT UPDATES?

if action selection is greedy, then no means for exploration.

If initialized at same Q for all S , then the two will make the same moves.

However, since initialization is hard to be random, the actions selected will be different & weight updates might be different too. & since greedy, then no shot of converging to same solution.