# Reinforcement Learning CS6700

Fall 2018

## Assignment 2

## Report

Author: Varun Sundar

EE16B068

24th September 2018

# 1    Implementation and Technical Notes

The code uses python 3.6 with sub-modules for questions. The repository adheres to the following:

- *Numpy style* documentation for the module and exposed functions.

- A *requirements.txt* for pip installing packages.

- Reproducible logs and reports.

- Code at github.com/varun19299/CS6700_Reinforcement_Learning

# 2    Question 1: Taxi Driver Problem

## 2.1   Part 1

Dynamic Programming via the compact Bellman operators was used to solve this problem.

We implement $T(J)$ by the following vectorised code:

```
np.amax(np.sum(r*P+P*np.expand_dims(J.T,2),axis=1),axis=1)
```

We get past the fact that at *town B* you cannot take action 3 (or our action 2) , by settting the rewards for action 2 (for B) as zero. Also note that since python indexing begins at zero, so do our numbering of states, stages and actions.

Where, $r$ is the reward matrix of shape (states, states, actions); $P$ is the probability matrix of shape (states, states, actions); $J$ is the set of states.

We do not assume the policy to be stationary (stage independent), however, this turns out to be the case in the optimal policy.

The results may be reproduced by running:

```
python3 q1.py --stages 10
```

## 2.2   Part 2

The optimal policy and rewards stage wise, for $N = 10$:

```
Starting with end stage costs as [ 0.   0.   0.]
Values of J [ 16.    15.     4.5] at stage 9
Optimal policy is at stage 9 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 26.25     29.40625  18.28125] at stage 8
Optimal policy is at stage 8 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 38.265625    43.51367188  29.453125  ] at stage 7
Optimal policy is at stage 7 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 49.859375    57.30688477  41.44128418] at stage 6
Optimal policy is at stage 6 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
```

```
Values of J [ 61.65032959   70.84981537   53.24645233] at stage 5
Optimal policy is at stage 5 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 73.44839096   84.17463732   65.14957047] at stage 4
Optimal policy is at stage 4 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 85.29898071   97.31518024   77.06275252] at stage 3
Optimal policy is at stage 3 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [  97.18086661   110.29839104    89.0057004 ] at stage 2
Optimal policy is at stage 2 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 109.09328351   123.1477526    100.96786822] at stage 1
Optimal policy is at stage 1 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Optimal policy is {'state 0': 'action 1', 'state 1': 'action 1', 'state
   2': 'action 2'}
```

Clearly, the optimal policy is stationary and equal to:

```
Optimal policy is {'state 0': 'action 1', 'state 1': 'action 1', 'state
   2': 'action 2'}
```

ie, ...,

- For towns A and B, Go to the nearest taxi stand and wait in line.

- For town C, Wait for a call from the dispatcher.

For $N = 20$;

```
Starting with end stage costs as [ 0.   0.   0.]
Values of J [ 16.    15.     4.5] at stage 19
Optimal policy is at stage 19 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 26.25     29.40625   18.28125] at stage 18
Optimal policy is at stage 18 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 38.265625     43.51367188   29.453125  ] at stage 17
Optimal policy is at stage 17 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 49.859375     57.30688477   41.44128418] at stage 16
Optimal policy is at stage 16 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 61.65032959   70.84981537   53.24645233] at stage 15
Optimal policy is at stage 15 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 73.44839096   84.17463732   65.14957047] at stage 14
Optimal policy is at stage 14 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 85.29898071   97.31518024   77.06275252] at stage 13
```

```
Optimal policy is at stage 13 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [  97.18086661  110.29839104   89.0057004 ] at stage 12
Optimal policy is at stage 12 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 109.09328351  123.1477526   100.96786822] at stage 11
Optimal policy is at stage 11 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 121.03057586  135.88310551  112.94817246] at stage 10
Optimal policy is at stage 10 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 132.98937416  148.52138909  124.94340833] at stage 9
Optimal policy is at stage 9 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 144.96639125  161.07701436  136.9515065 ] at stage 8
Optimal policy is at stage 8 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 156.95894887  173.56225617  148.9705143 ] at stage 7
Optimal policy is at stage 7 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 168.96473159  185.9875656   160.9988241 ] at stage 6
Optimal policy is at stage 6 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 180.98177784  198.36184213  173.03505106] at stage 5
Optimal policy is at stage 5 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 193.00841445  210.69266367  185.07802059] at stage 4
Optimal policy is at stage 4 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 205.04321752  222.9864829   197.12673118] at stage 3
Optimal policy is at stage 3 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 217.08497435  235.24879433  209.18033042] at stage 2
Optimal policy is at stage 2 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Values of J [ 229.13265238  247.48427659  221.23809236] at stage 1
Optimal policy is at stage 1 is {'state 0': 'action 1', 'state 1': '
   action 1', 'state 2': 'action 2'}
Optimal policy is {'state 0': 'action 1', 'state 1': 'action 1', 'state
   2': 'action 2'}
```

Again, the same policy.

Some comments on the rewards and policy:

- The rewards are unbounded - keep increasing - with stage. This is expected since there is no concept of termination in this problem.

- We donot end up with action 3 for *town B* - a sanity check.

- It is non-optimal to only go to the nearest taxi stand. We shall show this for a stage shortly.

## 2.3   Part 3

No, it is not optimal to to go to the nearest taxi stand, irrespective of the state.
    Assume this to be true. But,

$$J_1(C) = max_a E(r(C, j, a))$$
$$= max(2.5 + 0.5 + 0.25, 0.75 + 3 + 0.25, 3 + 1.5)$$
$$= max(3.25, 4, 4.5)$$

Hence, this is not so.

## 2.4   Code Blocks (pertinent only)

```
64  def  T(J,verbose=False,stage=None):
65          '''
66          Bellman  operator  for  maximising  reward
67
68          Args:
69          *  verbose:  Print  policy  each  time  T  is  operated.
70          *  stage:  Which  stage  to  operate  T  at.
71          '''
72          if  verbose  and  stage:
73              policy=np.argmax(np.sum(r*P+P*J.T[np.newaxis,:,np.newaxis],axis
                    =1),axis=1)
74              cost=np.amax(np.sum(r*P+P*J[np.newaxis,:,np.newaxis],axis=1),
                    axis=1)
75              print(f"Policy at stage {stage} is {policy}, cost at stage {cost
                    }")
76          else:
77              return np.amax(np.sum(r*P+P*J[np.newaxis,:,np.newaxis],axis=1),
                    axis=1)
78
79  def  read_optimal_policy(J_opt):
80          '''
81          Prints  policy  for  a  particular  optimal  J.
82
83          Agrs:
84          *  J_opt:  Optimal  reward.
85          '''
86          actions=  np.argmax(np.sum(r*P+P*J[np.newaxis,:,np.newaxis],axis=1),
                    axis=1)
87          return {f"state {state}":f"action {action}" for state,action in zip(
                    range(3),actions)}
```

# 3   Question 2

We now use similar code to question 1, except, including this into a class *Bellman* for further
flexibility.
    Modelling the grid-world:

- We use a 2D array, referenced by flattened indices for operating upon.

- Indices go from $(0,0)$ to $(9,9)$.

- Again, $P$, $r$, and $J$ are modelled by tensors.

- *Wormholes*: Correspond to probabilities of 1 towards transition.

- *Terminal*: Collect a one time, at transit reward of $+100$.

Results may be replicated by:

```
88    python3 q2.py --terminal 3 --stages -1 --supress 1
```

, similar for other parameters.

## 3.1   Code Blocks (pertinent only)

```
89  P=np.zeros((100,100,4))
90
91  for i in range(100):
92      # Up is +10
93      if (i<90):
94          P[i,i+10,0]=0.8
95      else:
96          P[i,i,0]=0.8
97      if i%10<9:
98          P[i,i+1,0]=0.1
99      else:
100         P[i,i,0]=0.1
101     if i%10>0:
102         P[i,i-1,0]=0.1
103     else:
104         P[i,i,0]=0.1
105
106     # Down is -10 ...
107     # Left is -1 ....
108     # Right is +1 ...
109
110 # Wormholes
111 for i in [32,42,52,62]:
112     for j in range(4):
113         P[i,i+1,j]=0
114         P[i,i-1,j]=0
115         P[i,i+10,j]=0
116         P[i,i-10,j]=0
117         P[i,0,j]=1
118
119         ....
120 # Terminal stage
```

```
121
122  P[terminal,terminal,:]=1
123  if (terminal%10)<9:
124      P[terminal,terminal+1,:]=0
125  if (terminal%10)>0:
126      P[terminal,terminal-1,:]=0
127  if (terminal<90):
128      P[terminal,terminal+10,:]=0
129  if (terminal>9):
130      P[terminal,terminal-10,:]=0
131
132  r=-1*np.ones((100,100,4))
133  r[:,terminal,:]=100
134  # Collect reward only once
135  r[terminal,terminal,:]=0
```

## 3.2 Part 1

We stop when the maximal absolute difference (np.max(np.abs())) of $J_i$ and $J_{i+1}$ falls below a certain $\epsilon$. For the sake of these three questions, we use $\epsilon = 1e - 6$.

The intuition for this follows from the fact that $T$ is a contraction mapping, hence its repeated operations produce a Cauchy sequence in metric complete space. We have used to strength this fact.

## 3.3 Part 2 : Code Snippet

Plots on pages 10 and 11.

```
136  def plot_convergence(J_array, terminal= args.terminal, stage=0,
          save_path=None,supress=False):
137      '''
138      Plot rewards, at stages.
139      Stages are inverted here for convinience, but nonetheless holds.
140      '''
141      J_array=np.array(J_array)
142      J_diff=np.max(np.abs(J_array[1:]-J_array[:-1]),axis=1)
143      iters=np.arange(1,len(J_diff)+1)
144      plt.plot(iters,J_diff)
145
146      plt.grid()
147
148      for j in range(len(J_diff)):
149          if (j<10 and j%3==0) or (j%10==0):
150              plt.text(j+1.15,J_diff[j]+0.15,s=f'value {J_diff[j]:.2f}')
151
152      plt.title("$max_s |J_{i+1}(s)    J_i(s)|$ vs iterations.")
153
154      # Save plot
155      if save_path:
```

```
156          plt . savefig ( os . path . join ( save_path , f " convergence - till -{ stage }.
                 png " ) )
157
158     if not supress :
159          plt . show ()
160     else :
161          plt . close ()
```

## 3.4   Part 3,4

We show the plots of $J$ and *actions* as heatmaps and quiver plots respectively.

Plots on pages 12 to 17.

Comments on the plots after absolute difference convergence below a pre-determined threshold:

- Wormholes are skipped wherever they lead away from the terminal state. (1 in case of (9,9), 2 in case of (3,0)).

- Reward to go is minimal at terminal state. (since once acquired, you terminate the game).

- The general policy is to either choose an apt wormhole or the direct shortest path to the terminal state.

- Since the probability in the intended action (Up when chosen up) is dominant, we see similar actions. If this were not the case, we could see non-obvious actions.

- Colliding into the walls (and thereby retaining state) is discouraged, unless you happen to be at the terminal state.

- Thus, this policy is "greedy" and does not take into account any time-variant phenomena, memory etc. Neither does solving this given grid require this.

## 3.5   Code Snippet for Part 3

```
162 def quiver_actions ( actions , terminal = args . terminal , stage =0 , save_path =
        None , supress = False ) :
163      ' ' '
164      Plot a quiver plot of the policy
165      ' ' '
166      def _action_u ( u ) :
167           ' ' '
168           Horz quiver
169           -1 ,1 if u == left or right
170           0 else
171           ' ' '
172           if u ==2:
173                return -1
```

```
174            elif u==3:
175                return 1
176            else:
177                return 0
178        def _action_v(u):
179            '''
180            Vert quiver
181            -1,1 if u == down or up
182            0 else
183            '''
184            if u==0:
185                return 1
186            elif u==1:
187                return -1
188            else:
189                return 0
190
191        X=Y=np.arange(0.5,10.5,1)
192        U= np.array([_action_u(a) for a in actions]).reshape((10,10))
193        V= np.array([_action_v(a) for a in actions]).reshape((10,10))
194        q=plt.quiver(X,Y,U,V)
195        plt.quiverkey(q,X=8, Y=8, U=1,label='Quiver key, length = 1',
                labelpos='E')
196        plt.title(f"Quiver state plot at stage {stage}")
197
198        major_ticks = np.arange(0, 10, 1)
199
200        # Wormholes 1
201        for j in range(3,8):
202            plt.scatter(2.5,j+0.5,s=225,color=colors['red'])
203            plt.text(2.5, j + 0.5, 'IN1')
204        # Exit 1
205        plt.scatter(0.5,0.5,s=225,color=colors['maroon'])
206        plt.text(0.5, 0.5, 'OUT1')
207
208        # Wormholes 2
209        plt.scatter(7.5,1.5,s=225,color=colors['grey'])
210        plt.text(7.5, 1.5, 'IN2')
211        plt.scatter(7.5,9.5,s=225,color=colors['lightgrey'])
212        plt.text(7.5, 9.5, 'OUT2')
213
214        # Terminal State
215        a=args.terminal//10+0.5
216        b=args.terminal%10+0.5
217        plt.scatter(b,a,s=256,color=colors['green'])
218        plt.text(b,a, 'TERMINAL')
219
220        plt.xlim((0,10))
221        plt.ylim((0,10))
```

```
222     plt.xticks(major_ticks)
223     plt.yticks(major_ticks)
224
225     plt.grid(True)
226
227     if save_path:
228         plt.savefig(os.path.join(save_path,f"quiver-{stage}.png"))
229     if not supress:
230         plt.show()
231     else:
232         plt.close()
```

# 4    References

- Classroom lectures.

- Bertsekas: *Dynamic Programming and Optimal Control, Vol 2, 3rd ed.*

- Numpy, Matplotlib Dev Documentation.

Figure 1: Convergence of $J_i$ till $N = 10$ for **Terminal State (9,9)**



Figure 2: Convergence of $J_i$ till $N = 25$ for **Terminal State (9,9)**

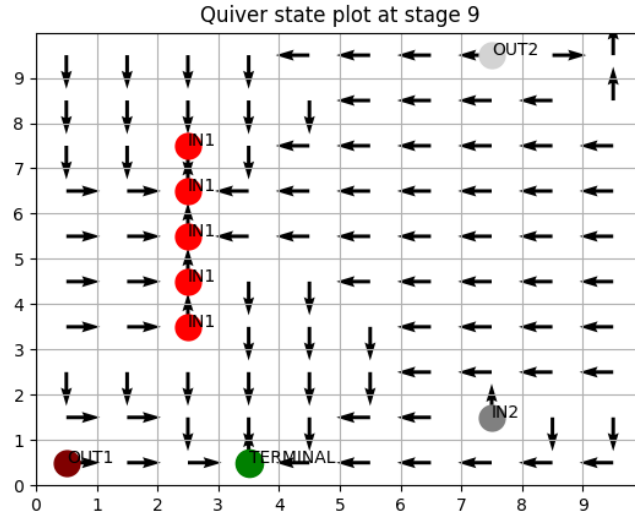Figure 4: Convergence of $J_i$ till $N = 10$ for **Terminal State (3,0)**



Figure 5: Convergence of $J_i$ till $N = 25$ for **Terminal State (3,0)**
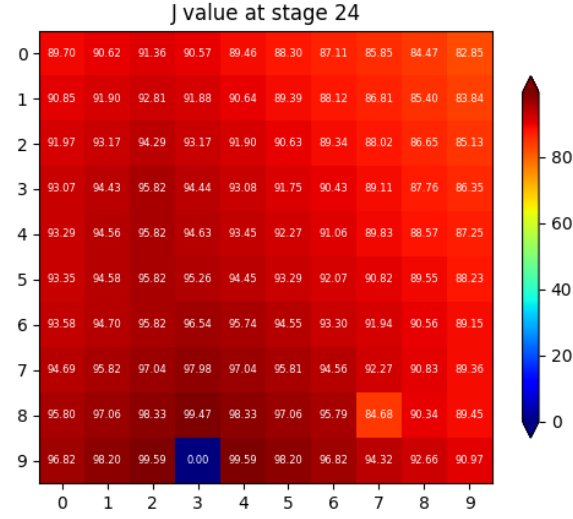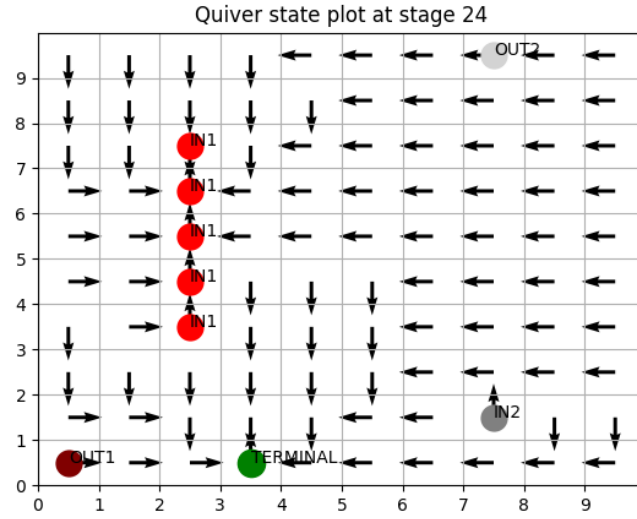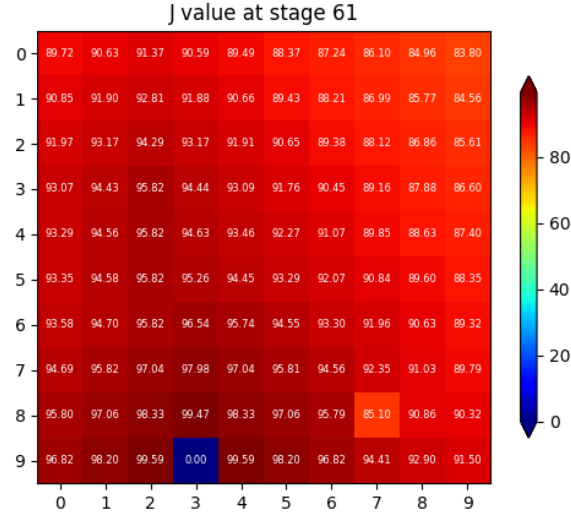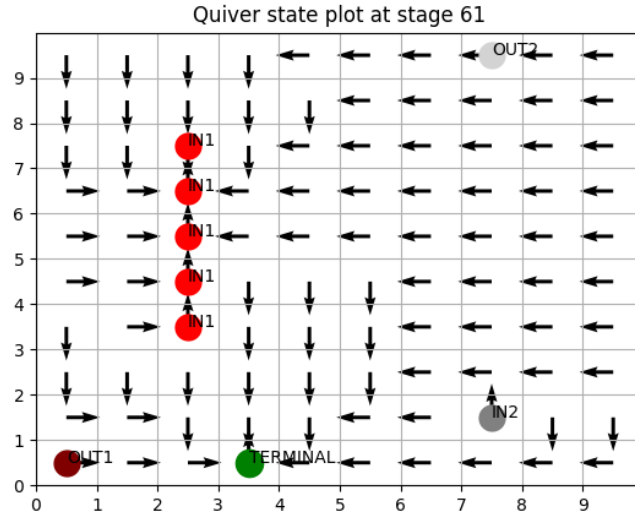
Figure 7: Heat-map of $J_i$ till $N = 10$ for **Terminal State (9,9)**



Figure 8: Quiver plot of $\pi$ till $N = 10$ for **Terminal State (9,9)**

Figure 9: Heat-map of $J_i$ till $N = 25$ for **Terminal State (9,9)**



Figure 10: Quiver plot of $\pi$ till $N = 25$ for **Terminal State (9,9)**

Figure 11: Heat-map of $J_i$ till absolute difference convergence for **Terminal State (9,9)**



Figure 12: Quiver plot of $\pi$ till absolute difference convergence for **Terminal State (9,9)**

Figure 13: Heat-map of $J_i$ till $N = 10$ for **Terminal State (3,0)**



Figure 14: Quiver plot of $\pi$ till $N = 10$ for **Terminal State (3,0)**

Figure 15: Heat-map of $J_i$ till $N = 25$ for **Terminal State (3,0)**



Figure 16: Quiver plot of $\pi$ till $N = 25$ for **Terminal State (3,0)**

Figure 17: Heat-map of $J_i$ till absolute difference convergence for **Terminal State (3,0)**



Figure 18: Quiver plot of $\pi$ till absolute difference convergence for **Terminal State (3,0)**