

MATH4406 HW4

Adam James Murray, 42625364

September 26th, 2014

1. Consider a two front war. You have a finite force and need to distribute it between the two fronts. The opposing forces have almost impenetrable strongholds in their capitals, however their leadership is not imperialist in nature, and will not in general pursue you unless you invade their homeland. Your goal is to maximise territory and minimise casualties. We can conflate these two variables, since a loss in men will more than likely correspond to a loss in territory, likewise a gain in territory will add the forces of the captured lands to your army. Your actions are to divert numbers of forces between the fronts. The strongholds at the heart of enemy nations will see your army suffer high casualties and be chased back across the empire. Our discounting factor here is attrition - as the war rages on, the size of your and your opponent's fighting forces dwindles, as does their morale. For small armies ($\lambda \approx 0$), the cost of pushing to the enemy capital is low and the territory gain is high, hence the optimal policy is to push forward at all costs and storm the capital. For large armies ($\lambda \approx 1$), the casualties are high and the territory gain low, and the optimal policy is to perhaps hedge your bets and only practice your imperialism on the local peasants.

2. The first observation to make is that given the symmetry of the problem, we would expect an optimal policy to be symmetric about the origin (state 0) since the cost of the states on the left is exactly mirrored by the costs of the states on the right. For example, if the optimal policy dictated that we choose action -1 in state -4 , then we would expect the optimal action to be 1 in state 4 .

Furthermore, when the discounting factor is close to zero, the high costs at the boundary points are less of a factor in our decision, and hitting the adjacent states with probability 0.5 will most likely be the best we can do. Hence in this case we would expect an optimal policy to drive the system to the boundary.

For a λ close to 1 , the penalty for hitting the boundary is more severe and in this case we would expect the opposite, i.e. that the system will try to avoid the boundary, and make do with the small rewards (or costs) associated with staying near the origin.

3-5. See appendix for code used.

7. Running each algorithm for $\lambda \in (0, 1)$ in increments of 0.1, we found that each method agreed on the optimal policies. We present the optimal policies in the following table.

λ	Optimal Policy
0.01 - 0.63	3, -1, -1, -1, -1, 0, 1, 1, 1, 1, -3
0.64 - 0.85	3, 1, -1, -1, -1, 0, 1, 1, 1, -1, -3
0.86 - 0.92	3, 1, 1, -1, -1, 0, 1, 1, -1, -1, -3
0.93 - 0.95	3, 1, 1, 1, -1, 0, 1, -1, -1, -1, -3
0.96 - 0.99	3, 1, 1, 1, 1, 0, -1, -1, -1, -1, -3

When running the policy evaluation we used the parameters $\varepsilon = 0.001$ with a break condition after 1000 iterations. We have also denoted our decision at the origin as 0 to represent that the choice made will not effect the overall behaviour of the policy.

Inspecting the optimal policies above, we see that for low λ , i.e. $\lambda \in (0.01, 0.63)$, the optimal policy is simply to drive the system to the end points. This is expected since the discounting factor nullifies the effect of the high cost at the end points.

For very high λ , i.e. $\lambda \in (0.96, 0.99)$, we find that the optimal policy is to try and keep the system as close as possible to the origin, as we would expect, given that in these cases the extremely high cost of the boundary is still present.

We then observe that in the interim cases between the two extreme, there is a systematic development of optimal policies, with the optimal deviation from the origin inversely proportional to the discounting factor λ .

Further note that all the optimal policies are symmetric about the origin. This is to be expected given the symmetries of the model.

Appendix

./MATH4406/HW4.py

```
1 import string
2 import numpy as np
3 import copy
4 np.set_printoptions(suppress=1,precision=5)
5
6 def GenerateDecisionRule(index, stateSpace):
7     #make sure our index is sensible
8     index = index%(pow(2,len(stateSpace))-2)
9     #convert the index to binary string with appropriate padding
10    rule = '{:0{pad}b}'.format(index, pad=(len(stateSpace)-2))
11    #convert rule string to integer list and convert zeros to -1s
12    rule = [(string.atoi(rule[j]) - (1-string.atoi(rule[j]))) for j in xrange(len(rule))]
13    #add in the boundary actions
14    rule.append(-3)
15    rule.insert(0,3)
16    return rule
17
18 def TransProb(state, newState, action, stateSpace):
19    if state == 0 and (newState == 1 or newState == state):
20        return 0.5
21    elif state == len(stateSpace)-1 and (newState == len(stateSpace)-2 or newState == state):
22        return 0.5
23    elif newState == state-1 and action == -1:
24        return 0.75
25    elif newState == state+1 and action == -1:
26        return 0.25
27    elif newState == state-1 and action == 1:
28        return 0.25
29    elif newState == state+1 and action == 1:
30        return 0.75
31    else:
32        return 0
33
34 def GenerateProbabilityMatrix(decisionRule, stateSpace):
35    #create an iterable for the dimensions of our matrix
36    dim = xrange(len(stateSpace))
37    #initialise the matrix
38    probMatrix = [[0 for j in dim] for i in dim]
39    #populate matrix withh transition probabilities
40    for i in dim:
41        for j in dim:
42            probMatrix[i][j] = TransProb(i, j, decisionRule[i], stateSpace)
43    return np.matrix(probMatrix)
44
45 def GenerateRewardVector(decisionRule, stateSpace):
46    rewardVector = [decisionRule[i]*stateSpace[i] for i in xrange(len(stateSpace))]
47    return np.matrix(rewardVector).T
48
49 def CalculateDiscountedReward(lmda, state, action, vOld, stateSpace):
50    total = 0
51    for i in xrange(len(stateSpace)):
52        total = total + lmda*TransProb(state, i, action, stateSpace)*vOld[i]
53    return total
54
55
56 #BRUTE FORCE METHOD
57 def BruteForceSolver(lmda, stateSpace):
58    S = stateSpace
59    #start our search at negative infinity
60    highestValue = np.NINF
61    finalRule = 0
62
63    for i in xrange(pow(2, len(S))-2):
64        rule = GenerateDecisionRule(i, S)
65        P = GenerateProbabilityMatrix(rule, S)
66        r = GenerateRewardVector(rule, S)
67        A = np.identity(len(r)) - lmda*P
68        Ainv = np.linalg.inv(A)
69        value = sum(Ainv*r)
70        if value >= highestValue:
71            highestValue = value
72            #set the middle action as zero for readability
73            rule[(len(S)-1)/2] = 0
74            finalRule = rule
75            #print rule
76
77    return finalRule
78
79
80 #VALUE ITERATION ALGORITHM
81 def RunPolicyEvaluationAlg(lmda, epsilon, stateSpace, maxItr):
82    S = stateSpace
83    vOld = [0 for i in xrange(len(S))]
84    vNew = [0 for i in xrange(len(S))]
85    policy = [0 for i in xrange(len(S))]
86    policy[0] = 3
```

```

89     policy[len(S)-1] = -3
91     #loop a max of maxItr times
92     for i in xrange(maxItr):
93         for j in xrange(len(S)):
94             #boundary cases
95             if j == 0:
96                 vNew[j] = S[j]*3 + CalculateDiscountedReward(lmda, j, 3, vOld, S)
97             elif j == len(S)-1:
98                 vNew[j] = S[j]*(-3) + CalculateDiscountedReward(lmda, j, 3, vOld, S)
99             else: #all other cases
100                 a1 = S[j]*(1) + CalculateDiscountedReward(lmda, j, 1, vOld, S)
101                 a2 = S[j]*(-1) + CalculateDiscountedReward(lmda, j, -1, vOld, S)
102                 if a1 > a2:
103                     vNew[j] = a1
104                     policy[j] = 1
105                 else:
106                     vNew[j] = a2
107                     policy[j] = -1
108             if np.linalg.norm(np.matrix(vNew) - np.matrix(vOld)) < (epsilon*(1-lmda))/(2*lmda):
109                 break
110             vOld = copy.copy(vNew)
111
112     policy[(len(S)-1)/2] = 0
113     return policy
114
115
116 #POLICY ITERATION ALGORITHM
117 def RunPolicyIterationAlg(lmda, stateSpace):
118     #choose the start policy
119     policyOld = GenerateDecisionRule(0, stateSpace)
120     policyNew = GenerateDecisionRule(1, stateSpace)
121     while policyOld != policyNew:
122         policyOld = copy.copy(policyNew)
123         P = GenerateProbabilityMatrix(policyOld, stateSpace)
124         r = GenerateRewardVector(policyOld, stateSpace)
125         A = np.identity(len(r)) - lmda*P
126         Ainv = np.linalg.inv(A)
127         v = Ainv*r
128         vOld = np.NINF
129         for i in xrange(pow(2, len(stateSpace)-2)):
130             policy = GenerateDecisionRule(i, stateSpace)
131             P = GenerateProbabilityMatrix(policy, stateSpace)
132             r = GenerateRewardVector(policy, stateSpace)
133             vNew = sum(r + lmda*P*v)
134             if vNew > vOld:
135                 policyNew = GenerateDecisionRule(i, stateSpace)
136                 vOld = vNew
137             #print policyOld
138             #print policyNew
139     policyNew[(len(stateSpace)-1)/2] = 0
140     return policyNew
141
142
143 S = range(-5,6)
144 epsilon = 0.001
145 maxItr = 1000
146
147 print '----\n Solutions for brute force evaluation'
148 lmda = 0.01
149 while lmda < 1:
150     print "{:.2f}".format(lmda), BruteForceSolver(lmda, S)
151     lmda = lmda + 0.01
152
153 print '----\n Solutions for policy evaluation'
154 lmda = 0.01
155 while lmda < 1:
156     print "{:.2f}".format(lmda), RunPolicyEvaluationAlg(lmda, epsilon, S, maxItr)
157     lmda = lmda + 0.01
158
159 print '----\n Solutions for policy iteration'
160 lmda = 0.01
161 while lmda < 1:
162     print "{:.2f}".format(lmda), RunPolicyIterationAlg(lmda, S)
163     lmda = lmda + 0.01

```

```

output:
----
2     Solutions for brute force evaluation
4     0.01 [3, -1, -1, -1, -1, 0, 1, 1, 1, 1, -3]
5     0.02 [3, -1, -1, -1, -1, 0, 1, 1, 1, 1, -3]
6     0.03 [3, -1, -1, -1, -1, 0, 1, 1, 1, 1, -3]
7     0.04 [3, -1, -1, -1, -1, 0, 1, 1, 1, 1, -3]
8     0.05 [3, -1, -1, -1, -1, 0, 1, 1, 1, 1, -3]
9     0.06 [3, -1, -1, -1, -1, 0, 1, 1, 1, 1, -3]
10    0.07 [3, -1, -1, -1, -1, 0, 1, 1, 1, 1, -3]
11    0.08 [3, -1, -1, -1, -1, 0, 1, 1, 1, 1, -3]
12    0.09 [3, -1, -1, -1, -1, 0, 1, 1, 1, 1, -3]
13    0.10 [3, -1, -1, -1, -1, 0, 1, 1, 1, 1, -3]
14    0.11 [3, -1, -1, -1, -1, 0, 1, 1, 1, 1, -3]

```



```
300 || 0.95 [3, 1, 1, 1, -1, 0, 1, -1, -1, -1, -3]
      || 0.96 [3, 1, 1, 1, 1, 0, -1, -1, -1, -1, -3]
302 || 0.97 [3, 1, 1, 1, 1, 0, -1, -1, -1, -1, -3]
      || 0.98 [3, 1, 1, 1, 1, 0, -1, -1, -1, -1, -3]
304 || 0.99 [3, 1, 1, 1, 1, 0, -1, -1, -1, -1, -3]
```