

COMP 767 Reinforcement Learning

Assignment 1

Ali Rezagholizadeh

January 30, 2019

1 Question 1(b)

In this paper, authors tried to find problem-independent and problem-dependent boundary for the *finite time expected regret* of Thompson Sampling. we ignore describing Thompson Sampling method in this assignment because of being out of the scope of requirements in this question.

1.1 Summary of results

While several works had been done on finding boundary for *expected regret* in the bandit problem like *Lai and Robbins* (who proposed the lower bound of any bandit problem eq.1) [15]¹, *Auer et al.* (which proposed upper bound for UCB1 method eq.2) [4], and *Kaufmann et al.* [13 and 14], the authors of this paper achieved a problem-dependent upper boundary (eq.3) which approaches asymptotic lower bound of *Lai and Robbins* and match those provided by recent work for Thompson Sampling [14].

They also achieved a problem-independent upper bound (eq. 4) which approaches problem-independent lower bound of $\Omega(\sqrt{NT})$ which achieved by Bubeck et al. [5] for this problem.

$$E[R(T)] \geq \left[\sum_{i=2}^N \frac{\Delta_i}{d(\mu_i, \mu_1)} + o(1) \right] \ln T \quad (1)$$

$$E[R(T)] \leq \left[8 \sum_{i=2}^N \frac{1}{\Delta_i} \right] \ln T + (1 + \pi^2/3) \left(\sum_{i=2}^N \Delta_i \right) \quad (2)$$

$$E[R(T)] \leq (1 + \epsilon) \sum_{i=2}^N \frac{\ln T}{d(\mu_i, \mu_1)} \Delta_i + O\left(\frac{N}{\epsilon^2}\right) \quad (3)$$

$$E[R(T)] \leq O(\sqrt{NT \ln T}) \quad (4)$$

These two problem-dependent and problem-independent bounds (eq.3 and 4) were brought as Theorems in the paper.

1.2 Main steps and ideas in the proof

Authors constrained their work to have Thompson Sampling with Bernoulli reward (reward of 0 or 1), but as they confirmed this can be extended to general reward distribution using [1]. This limitation to have Bernoulli reward helps authors to get a benefit from conjugate of Bernoulli, Beta distribution, as prior distribution ($p(R)$) of the reward. This makes posterior ($p(R | h)$), where h is: $a_1, r_1, a_2, r_2, \dots, a_t, r_t$) also being Beta distribution. Posterior of each arm changes whenever the arm is played. Whatever the

¹This is the reference number used in the paper.

reward achieved by playing an arm be one (success) (or zero (fail)), then the mean of its posterior distribution (expected reward according to this current distribution) approaches to one (or zero) and it becomes tighter at the mean point.

Sampling from these distributions and then selecting the next arm to play according to the value of each sample, makes possible for us to estimate the expected regret (or at least find a boundary for it) in time T . The authors did this by getting benefit from these distributions.

In other side, the authors took beautiful strategy by considering two arbitrary points, x_i and y_i for each arm so that $\mu_i < x_i < y_i < \mu_1$, where μ_i is the value function of arm $i \neq 1$ (here, it should be stated that the authors, without loss of generality, assumed that the first arm is the unique optimal arm, i.e., $\mu^* = \mu_1 > \argmax_{i \neq 1} \mu_i$).

By considering of these x_i and y_i , they constructed the foundation of their inequalities used in the steps (lemmas) throughout the way of proving the Theorems.

The concept behind using x_i and y_i could be both:

1. Give the opportunity to design the boundary whatever they want in the share point:

$$E[K_i(T)] \leq \frac{24}{\Delta_i'^2} + \sum_{j=0}^{T-1} \Theta \left(e^{-\Delta_i'^2 j/2} + \frac{1}{(j+1)\Delta_i'^2} e^{-D_i j} + \frac{1}{e^{\Delta_i'^2 j/4} - 1} \right) + L_i(T) + 1 + \frac{1}{d(x_i, \mu_i)} + 1. \quad (5)$$

, where the way of obtaining problem-dependent bound and problem-independent bound are separating.

2. assure that after a time t (before T), so that *for every* $i \neq 1$, $K_i(t) > L_i(T)$, there exist such inequality in the aforementioned share point. This because $L_i(T)$ is the function of T , x_i , and y_i .

The main steps through proving theorems are four lemmas whose aim is to prove the share equation (eq. 5) and then by selecting suitable X_i and Y_i (settings in (6)) they reached to the mentioned theorems. It is worth mentioning that, through the way, author got advantages of Chernaff-Hoeffding bound and 'Beta-Binomial trick' facts a lot.

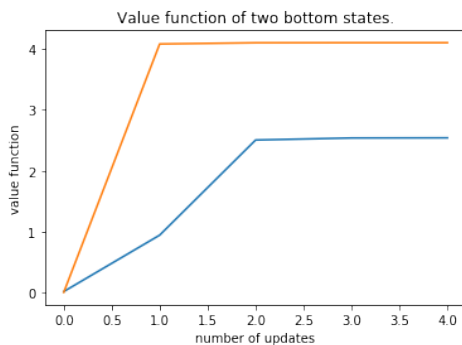
$$\begin{cases} d(x_i, \mu_1) = d(\mu_i, \mu_1)/(1 + \epsilon) \text{ and} \\ d(x_i, y_i) = d(x_i, \mu_1)/(1 + \epsilon) = d(\mu_i, \mu_1)/(1 + \epsilon)^2 \quad \text{to reach problem-dependent bound,} \\ x_i = \mu_i + \frac{\Delta_i}{3} \text{ and } y_i = \mu_1 - \frac{\Delta_i}{3} \quad \text{to reach problem-independent bound.} \end{cases} \quad (6)$$

2 Question 2(a)

I could implement Policy iteration completely. I bring the result of my code here and my raw code after:

1. $p = 0.7$ and $n = 5$: Total update in policy iteration: 4 Final greedy policy and its value-function:

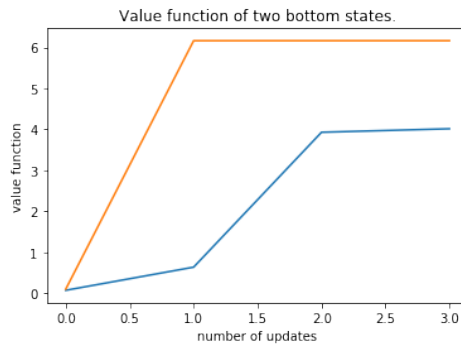
term	R	R	R	term	0	5.2405887925621055	6.875761363794614	8.723676358173222	0
R	R	R	R	U	4.156965591764831	5.478203156223082	6.477986299407242	7.549084423452747	8.724719871571663
R	R	R	U	U	3.8682767418740753	4.82266935665731	5.609188915350512	6.481219238604924	6.891213164547389
R	R	U	U	U	3.3199941309863172	4.1399674406401745	4.826129030029215	5.510483427105133	5.455393907238159
R	U	U	U	U	2.540626167075174	3.3214355869812446	3.886252443697519	4.382699323992172	4.100969454668087



2. $p = 0.9$ and $n = 5$:
Total update in policy iteration: 3 Final greedy policy:

term	R	R	R	term	0	7.128892329611076	8.349886293709046	9.630877789108	0
R	R	R	R	U	5.790348385629835	6.7827685259004475	7.634687771942954	8.578015472911629	9.63093775009695
R	R	R	U	U	5.275605122835685	6.0537192831497935	6.800231175069826	7.634869923265149	8.352551201646104
R	R	U	U	U	4.693923689708901	5.387312419175714	6.0539062034463065	6.78819738182139	7.244931180934402
R	U	U	U	U	4.01335752282207	4.694002915325792	5.278436309672564	5.908637001806977	6.164349540668548

3. $p = 0.7$ and $n = 50$:
It's run took a lot of time...



2.1 Code:

```

1 import numpy as np
2
3 def Setting_Rewards():
4     Rewards = [0, 1, 10]
5
6     return Rewards
7
8
9 def Setting_Actions():
10     Actions = ["U", "D", "L", "R"]
11
12     return Actions
13
14
15 def Convert_position_list_to_str(s):
16     # s is [,] position.
17     s_str = '['
18     s_str += ','.join(str(x) for x in s)
19     s_str += ']'
20
21     return s_str
22
23 def Convert_position_str_to_list(s_str):
24     r = int(s_str[1])
25     c = int(s_str[2])
26
27     s = [r, c]
28
29     return s
30
31 def greedy_policy_of_current_policy(S_b, S_t, V):
32     import random
33     Greedy_Policy = {}
34     for s in S_b:
35         Action_set = Deterministic_Actions(s, V)
36         if (Action_set == "Error"):
37             print("Error occurred in obtaining greedy policy in
state: ", s)

```

```

38         else:
39             s_str = Convert_position_list_to_str(s)
40             if(len(Action_set) > 0):
41                 n_a = len(Action_set)
42                 direction_n = int(random.uniform(0,n_a))
43                 direction = Action_set[direction_n]
44                 d = {s_str: direction}
45                 Greedy_Policy.update(d)
46
47     for s in S_t:
48         s_str = Convert_position_list_to_str(s)
49         d = {s_str: ""}
50         Greedy_Policy.update(d)
51
52
53
54     # Greedy_Policy is a dict like: {"[00]": 'U'}
55     return Greedy_Policy
56
57
58 def Determine_States(n):
59     # This is determined according the requirements of the ass1
60     # - q2(a) in Comp 767.
61     # the states are position like [[,],[,]...]
62     S_t = [[0,0] , [0,n-1]]
63     # S_t is terminal state.
64     S_b = [[i,j] for i in range(n) for j in range(n)]
65     S_b.pop(n-1)
66     S_b.pop(0)
67     # S_b is in between - not terminal state.
68
69     return S_t , S_b
70
71 def Inital_value_function(S_b,S_t):
72     # S_t is terminal states, and 'S_b' is between sates.
73     # Each is list of positions. [[,] [,] ..]
74     import math
75     n1 = len(S_b)
76     n2 = len(S_t)
77
78     n = n1+n2
79     width = int(math.sqrt(n))
80     Value_Function = []
81     for i in range(width):
82         r = [0 for j in range(width)]
83         Value_Function.append(r)
84
85     return Value_Function
86
87 def Deterministic_Actions(s,V):
88     # this is deterministic action selection according to policy
89     Pai.

```

```

89 # So, it select an action according to value function of
    that policy.
90 # input: s is a position of one point, list of two numbers (
    position in gride- which is starting from 0)      -
91 #         ... V is Value function of a policy (is matrix of
    gride size).
92 import numpy as np
93 r = s[0]
94 c = s[1]
95
96 All_logical_neighbors = [[r-1, c], [r+1, c], [r, c-1], [r,
    c+1]]
97 end_of_grid = len(V)
98 Possible_neighbors = [z for z in All_logical_neighbors if (z
    [0]>-1 and z[0]< end_of_grid and z[1]>-1 and z[1]<
    end_of_grid)]
99 max_neighbor_value = -10
100
101 for position in Possible_neighbors:
102     value = V[position[0]][position[1]]
103     if (max_neighbor_value < value):
104         max_neighbor_value = value
105
106 max_neighbor_pos = []
107 for position in Possible_neighbors:
108     value = V[position[0]][position[1]]
109     if (max_neighbor_value == value):
110         max_neighbor_pos.append(position)
111 # all neighbors who have the max value would enter to '
    max_neighbor_pos '.
112
113 action_set = ""
114 s_array = np.array(s)
115 for max_Neigh in max_neighbor_pos:
116     max_neighbor_arr = np.array(max_Neigh)
117     move = max_neighbor_arr - s_array
118
119     if (move[0] == -1):
120         action_set += "U"
121     elif (move[0] == 1):
122         action_set += "D"
123     else:
124         if (move[0] != 0):
125             print("Error in Choosing-deterministic_Action ,
    1.", move[0])
126             action_set = "Error"
127
128
129     if (move[1] == -1):
130         action_set += "L"
131     elif (move[1] == 1):
132         action_set += "R"
133     else:
134         if (move[1] != 0):

```

```

135         print("Error in Choosing_deterministic_Action ,
136         2." , move[1])
137         action_set = "Error"
138
139
140     return action_set
141
142
143 def Printing_ValueFunction(matrix_list):
144     # 'matrix_list' is [[ ,...][ ,...]....]
145     r_len = len(matrix_list)
146     c_len = len(matrix_list[0])
147     print("Value Function:" , end ='\n')
148     for i in range(r_len):
149         for j in range(c_len):
150             if(j==0):
151                 print(end=' | ')
152                 print(matrix_list[i][j] , end=' | ')
153             else:
154                 print(matrix_list[i][j] , end=' | ')
155         print(' ')
156     return
157
158 def Printing_Deterministic_Policy(g_p,S_t , S_b):
159     # g_p is dict of like {'[11]': "UDLR"}
160     import math
161     n1 = len(S_b)
162     n2 = len(S_t)
163
164     n = n1+n2
165     width = int(math.sqrt(n))
166
167     Greedy_Policy_mat = []
168     for i in range(width):
169         r = [ '' for j in range(width)]
170         Greedy_Policy_mat.append(r)
171
172     for s in S_b:
173         s_str = Convert_position_list_to_str(s)
174         Greedy_Policy_mat[s[0]][s[1]] = g_p[s_str]
175
176
177     for s in S_t:
178         #s_str = Convert_position_list_to_str(s)
179         Greedy_Policy_mat[s[0]][s[1]] = "term"
180
181
182     # printing
183     matrix_list = Greedy_Policy_mat
184     r_len = len(matrix_list)
185     c_len = len(matrix_list[0])
186
187     # 'value %3s - num of occurances = %d' % item

```



```

188     #'{0:10} ==> {1:10d}'.format(name, phone)
189
190
191     for i in range(r_len):
192         for j in range(c_len):
193             if (j==0):
194                 print(end=' | ')
195                 print("%5s |"% matrix_list[i][j], end='')
196             else:
197                 print("%5s |"% matrix_list[i][j], end='')
198             print(' ')
199
200
201     return
202
203
204
205 #

```

```

206 # DP part:
207 #
208
209
210 def Policy_Iteration(S_b, S_t, Rewards, Actions, Policy, V,
211                     Discount, P_Dynamic, e):
212     import math
213     Greedy_Policy = Policy
214     # Greedy_Policy is dict of like {'[12]':'U'}
215     Theta = 0.001
216
217     S_p = S_b + S_t
218     n = len(S_p)
219     n = int(math.sqrt(n))
220     policy_stable = False
221     # S_p is s+
222     Number_of_update_policy = 1
223     V_S_Left = []
224     V_S_Right = []
225     while(policy_stable == False):
226
227         print("===== update: ",
228               Number_of_update_policy, " =====")
229
230         print("Greedy policy of current policy: \n")
231         Printing_Deterministic_Policy(Greedy_Policy, S_t, S_b)
232
233         print("Computing V_i through the value function of
234               aforementioned greedy policy ... ..")
235         Delta = 2
236         while(Theta < Delta):
237             Delta = 0
238             for s in S_b:
239                 # s is [r,c] as the position of current state.

```

```

236         r_s = s[0]
237         c_s = s[1]
238         current_v = V[r_s][c_s]
239         s_str = Convert_position_list_to_str(s)
240         greedy_action = Greedy_Policy[s_str]
241
242         Greedy_Policy_ = E_greedy_policy(Greedy_Policy,
243     S_b, Actions, e)
244         actions, p_ = Greedy_Policy_[s_str]
245         if(greedy_action != ""):
246             # when s is not terminal
247             sum_ = 0
248             for a_i in range(len(Actions)):
249                 pi = p_[a_i]
250                 ac = Actions[a_i]
251                 sum_ = 0
252                 for s_ in S_p:
253                     for r in Rewards:
254                         p_dy = Dynamic_p(s_, r, s, ac,
255     P_Dynamic)
256                         if(p_dy != 0):
257                             #print("s, greedy_action, S_
258     , r, p",s, greedy_action, s_, r , p_dy)
259                             sum_ += p_dy*(r + Discount*
260     V[s_[0]][s_[1]])
261                             sum_ += pi * sum_
262
263                 new_v = sum_
264                 V[r_s][c_s] = new_v
265                 diff = abs(current_v - new_v)
266                 #print(" diff: ", diff, end=' - ')
267                 Delta = max(Delta, diff)
268
269             #print()
270             #print("Theta and Delta: ",Theta, Delta, end=' - ')
271             print("Value-function obtained from 'Iterative policy
272     evaluation' is: ")
273             Printing_ValueFunction(V)
274             V_S_Left.append(V[n-1][0])
275             V_S_Right.append(V[n-1][n-1])
276             print("_____")
277             print("Now updating greedy policy using this
278     Value_function achieved ... ..")
279             policy_stable = True
280             for s in S_b:
281                 s_str = Convert_position_list_to_str(s)
282                 current_action = Greedy_Policy[s_str]
283
284                 max_action_value = 0
285                 max_q_action = ""
286                 for a in Actions:
287                     sum_ = 0
288                     for s_ in S_p:
289                         for r in Rewards:
290                             p_dy = Dynamic_p(s_, r, s, a, P_Dynamic)

```

```

284         if(p_dy != 0):
285             sum_ += p_dy*(r + Discount* V[s_
[0]][s_[1]])
286
287             if(max_action_value <= sum_):
288                 max_action_value = sum_
289                 max_q_action = a
290
291
292         Greedy_Policy[s_str] = max_q_action
293
294         if(current_action != max_q_action):
295             policy_stable = False
296
297         print("Greedy policy of updated policy: \n")
298         Printing_Deterministic_Policy(Greedy_Policy, S_t, S_b)
299         if(policy_stable == False):
300             Number_of_update_policy += 1
301             print("Upadted policy is not stable so we go through
obtaining its Value-fuction, then updating g-policy again. "
)
302             print("
")
303             print("||||||||||||||||||||||||||||||||||||||||")
304             else:
305                 print("No update was needed. The last one is the
optimal policy.")
306
307
308
309
310         return V_S_Left, V_S_Right, Number_of_update_policy
311
312
313 #

```

```

314 def E_greedy_policy(Greedy_Policy, S_b, Actions, e):
315     Greedy_Policy_ = {}
316     for s in S_b:
317         s_str = Convert_position_list_to_str(s)
318         greedy_action = Greedy_Policy[s_str]
319         indx = 0
320         index = -1
321         for ac in Actions:
322             if(ac == greedy_action):
323                 index = indx
324                 indx += 1
325         p = [0 for i in range(len(Actions))]
326         p[index] = e
327         e_ = (1-e)/4
328         p_ = [i + e_ for i in p]
329         r = {s_str: [Actions, p_]}

```

```

330         Greedy_Policy_.update(r)
331
332     return Greedy_Policy_
333 def E_greedy(Greedy_Policy, S_b, Actions, e):
334     import random
335
336     for s in S_b:
337         s_str = Convert_position_list_to_str(s)
338         greedy_action = Greedy_Policy[s_str]
339
340         t = random.uniform(0,1)
341         if(t<e):
342             selection = greedy_action
343         else:
344             t_2 = random.uniform(0,1)
345             if(t_2 < 0.25):
346                 selection = Actions[0]
347             elif(t_2 < 0.5):
348                 selection = Actions[1]
349             elif(t_2 < 0.75):
350                 selection = Actions[2]
351             else:
352                 selection = Actions[3]
353
354     return selection
355
356 #

```

```

357 # Regarding Dynamic:
358 def Possible_Action_for_states(S_b,S_t):
359     import numpy as np
360     Possible_Actions_for_states = {}
361     terminal_1_str = Convert_position_list_to_str(S_t[0])
362     terminal_2_str = Convert_position_list_to_str(S_t[1])
363     for s in S_b:
364         r = s[0]
365         c = s[1]
366
367         All_logical_neighbors = [[r-1, c], [r+1, c], [r, c-1],
368 [r, c+1]]
369         end_of_grid = len(V)
370         Possible_neighbors = [z for z in All_logical_neighbors
371 if(z[0]>-1 and z[0]< end_of_grid and z[1]>-1 and z[1]<
372 end_of_grid)]
373         max_neighbor_value = -10
374
375         s_array = np.array(s)
376         action_set = ""
377         Next_state_set = []
378         Rewards = []
379         s_array = np.array(s)
380         for pos_N in Possible_neighbors:
381             pos_neighbor_arr = np.array(pos_N)

```

```

379         move = pos_neighbor_arr - s_array
380
381         if(move[0] == -1):
382             action_set += "U"
383             N_state_str = Convert_position_list_to_str(pos_N
384     )
385             Next_state_set.append(N_state_str)
386             if(N_state_str == terminal_1_str):
387                 Rewards.append(1)
388             elif(N_state_str == terminal_2_str):
389                 Rewards.append(10)
390             else:
391                 Rewards.append(0)
392         elif(move[0] == 1):
393             action_set += "D"
394             N_state_str = Convert_position_list_to_str(pos_N
395     )
396             Next_state_set.append(N_state_str)
397             if(N_state_str == terminal_1_str):
398                 Rewards.append(1)
399             elif(N_state_str == terminal_2_str):
400                 Rewards.append(10)
401             else:
402                 Rewards.append(0)
403         else:
404             if(move[0] != 0):
405                 print("Error in
406     Choosing_deterministic_Action , 1." , move[0])
407                 action_set = "Error"
408
409
410         if(move[1] == -1):
411             action_set += "L"
412             N_state_str = Convert_position_list_to_str(pos_N
413     )
414             Next_state_set.append(N_state_str)
415             if(N_state_str == terminal_1_str):
416                 Rewards.append(1)
417             elif(N_state_str == terminal_2_str):
418                 Rewards.append(10)
419             else:
420                 Rewards.append(0)
421         elif(move[1] == 1):
422             action_set += "R"
423             N_state_str = Convert_position_list_to_str(pos_N
424     )
425             Next_state_set.append(N_state_str)
426             if(N_state_str == terminal_1_str):
427                 Rewards.append(1)
428             elif(N_state_str == terminal_2_str):
429                 Rewards.append(10)
430             else:
431                 Rewards.append(0)
432         else:

```

```

428         if(move[1] != 0):
429             print("Error in
Choosing_deterministic_Action , 2." , move[1])
430             action_set = "Error"
431
432         s_str = Convert_position_list_to_str(s)
433         state_action = {s_str:[action_set , Next_state_set ,
Rewards]}
434         Possible_Actions_for_states.update(state_action)
435     for s in S_t:
436         action_set = ""
437         Next_state_set = []
438         Rewards = []
439         s_str = Convert_position_list_to_str(s)
440         state_action = {s_str:[action_set , Next_state_set ,
Rewards]}
441         Possible_Actions_for_states.update(state_action)
442     return Possible_Actions_for_states
443
444 def Dynamic_p(s_,r,s,a , P):
445     # s is like [1,1] ; a is like "U" ; r is like 0 ; s_ is like
[3,3]
446     # P is like: { '[01]':[Actions , Next_state , Rewards]}
447     # S_b
448     s_str = Convert_position_list_to_str(s)
449     Actions , Next_state , Rewards = P[s_str]
450     #print(Actions)
451     indx = 0
452     index = -1
453     for ac in Actions:
454         if(ac == a):
455             index = indx
456             indx += 1
457
458     # now index is corresponding action to take , next state ,
next reward considering this action.
459     if(index == -1):
460         return 0
461     s__str = Convert_position_list_to_str(s_)
462     if(Next_state[index] != s__str):
463         return 0
464     if(Rewards[index] != r):
465         return 0
466
467     return 1
468
469     # for each 'current state' , 'action' and 'reward' , 'new
state' it should return the dynamic of
470 #

```

```

471
472 def Value_Iteration():
473

```

```

474     return
475
476 def Modified_Policy_Iteration() :
477
478     return
479
480
481 def Main() :
482     S_t, S_b = Determine_States(50)
483     Rewards = Setting_Rewards()
484     Actions = Setting_Actions()
485     V = Initial_value_function(S_b, S_t)
486     e = 0.7
487     Discount = 0.9
488     g_Policy = greedy_policy_of_current_policy(S_b, S_t, V)
489     P = Possible_Action_for_states(S_b, S_t) # this is created
490     V_S_Left, V_S_Right, Number_of_update_policy =
491     Policy_Iteration(S_b, S_t, Rewards, Actions, g_Policy, V,
492     Discount, P, e)
493
494     import matplotlib.pyplot as plt
495
496     # x axis values
497     x = [i for i in range(Number_of_update_policy)]
498     # corresponding y axis values
499     y_1 = V_S_Left
500     # plotting the points
501     plt.plot(x, y_1, label = "Botton left")
502
503     y_2 = V_S_Right
504     plt.plot(x, y_2, label = "Botton Right")
505     # naming the x axis
506     plt.xlabel('number of updates')
507     # naming the y axis
508     plt.ylabel('value function')
509
510     # giving a title to my graph
511     plt.title('Value function of two bottom states.')
512
513     # function to show the plot
514     plt.show()
515
516     return

```

Listing 1: My code