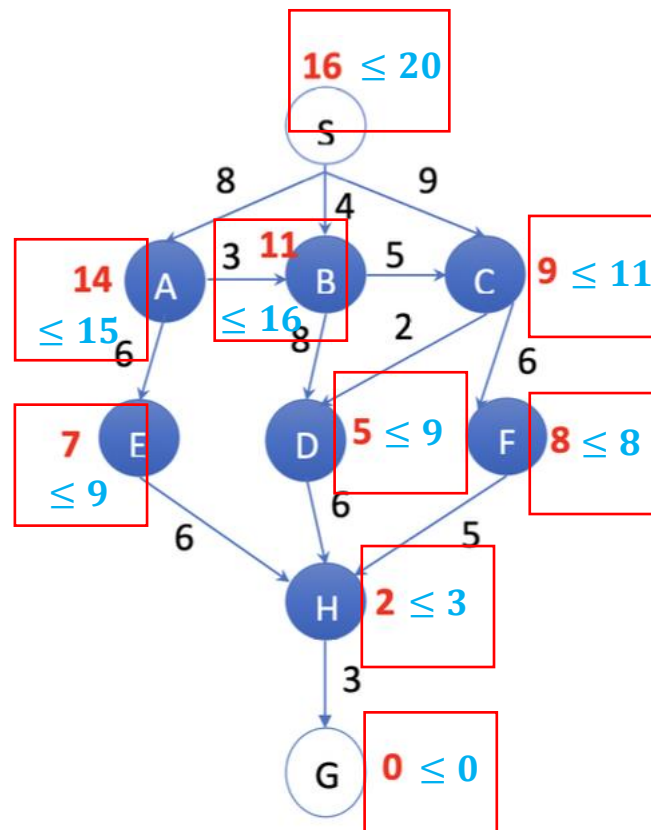


# CS420 ASSIGNMENT 2 SOLUTIONS

## QUESTION 1



(a) Yes, the heuristic is admissible.

Justification:

- For a heuristic to be admissible, the following must be satisfied:  

$$h(n) \leq TC^*(n, g)$$
- Since all the minimum cost to reach goal state from the current state (shown in light blue) is less than corresponding heuristic value (shown in red) of that state, the heuristic used is said to be admissible.

Node $n$	Minimum cost to reach goal from $n$ , $TC^*(n, g)$	$h(n)$
S	20	16
A	15	14
B	16	11
C	11	9
E	9	7
D	9	5

F	8	8
H	3	2
G	0	0

(b) No, the heuristic is not consistent.

**Justification:**

- For a heuristic to be consistent, the following must be satisfied:

$$h(n) \leq c(n,p) + h(p)$$

- That is, the heuristic value to reach goal state, g from the current state, n is no greater than the step cost of getting to state p plus the heuristic cost of reaching g from p.

Node $n$	$c(n,p) + h(p)$	$h(n)$
S	<ul style="list-style-type: none"> <li>- <math>c(S,A) + h(A) = 8 + 14 = 22</math></li> <li>- <math>c(S,B) + h(B) = 4 + 11 = 15 (&lt;16)</math></li> <li>...</li> </ul>	16

- Since the consistent heuristic must holds for every state and there is a violation, we conclude that the heuristic is not consistent.

(c) Part 1 : DFS Algorithm

- Final Solution Path :  $S \rightarrow A \rightarrow E \rightarrow H \rightarrow G$
- Cost of Solution Path :  $8 + 6 + 6 + 3 = 23$
- Open list's data structure : **Stack**

Step #	Open List (Stack)	Pop	Nodes to Add
1	$S$	$S$	$C^S, B^S, A^S$
2	$A^S, B^S, C^S$	$A^S$	$E^A$
3	$E^A, B^S, C^S$	$E^A$	$H^E$
4	$H^E, B^S, C^S$	$H^E$	$G^H$
5	$G^H, B^S, C^S$	$G^H$	Termination

Part 2 : A\* Algorithm

- Final Solution Path :  $S \rightarrow C \rightarrow D \rightarrow H \rightarrow G$
- Cost of Solution Path :  $9 + 2 + 6 + 3 = 20$
- Open list's data structure : **Priority Queue**

Step #	Open List (Priority Queue)	Dequeue	Nodes to Add
1	$S_{16}$	$S_{16}$	$A_{22}^S, B_{15}^S, C_{18}^S$
2	$B_{15}^S, C_{18}^S, A_{22}^S$	$B_{15}^S$	$D_{17}^B, C_{18}^B$ * Do not replace C in open list
3	$D_{17}^B, C_{18}^S, A_{22}^S$	$D_{17}^B$	$H_{20}^D$
4	$C_{18}^S, H_{20}^D, A_{22}^S$	$C_{18}^S$	$D_{16}^C, F_{23}^C$ * Re-open D in close list
5	$D_{16}^C, H_{20}^D, A_{22}^S, F_{23}^C$	$D_{16}^C$	$H_{19}^D$ * Replace H in open list
6	$H_{19}^D, A_{22}^S, F_{23}^C$	$H_{19}^D$	$G_{20}^H$
7	$G_{20}^H, A_{22}^S, F_{23}^C$	$G_{20}^H$	Terminated

- The priority of each state, n is based on the f-cost,  $f(n)$  which is defined as:

$$f(n) = g(n) + h(n)$$

That is, *sum of  $g(n)$ , cost from start state,  $s_0$  to current state,  $n$ , and  $h(n)$ , heuristic from current state,  $n$  to goal state,  $g$ .*

- Note: Since the heuristic is not consistent, re-opening of closed nodes is possible.

## QUESTION 2

### Clarification of term(s):

- 'First' refers to the integer at the front(leftmost element) of a list.

### Part 1

#### I. State representation

##### Describe specifically what is a state in this problem:

- Given N discs, a state could be any legal combination (a bigger disc should never be put on top of a smaller disc) of discs on top of each of the three towers.

##### How you would store it in a computer using a data structure, and justify the correctness of your state representation:

- One of the possible data structures for this state representation would be using three lists, storing the information of which tower is having which discs (of specific size) being stacked appropriately on it. Since we have N discs of different sizes, each can be assigned with an integer value from 1,2...,N corresponding to smallest to biggest disc. This allows us to represent each tower as an ordered list of integers. For example, the start state,  $s_0$  would be ([1,2...,N], [], []).

#### II. Actions

- Here, an action  $X_y$ , can be defined as **popping** the first integer(i.e., disc) from a list, X (i.e., tower) and **pushing** it to the front of another list, Y.
- An exhausted list of all possible actions for this search problem would be:

$$\text{Operators, } A = \{A_B, A_C, B_A, B_C, C_A, C_B\}$$

- For an action to be valid at a given state (not all actions are valid at a given each state):
  - The integer value being pushed to the front of another list must be strictly smaller than the integer currently located at the front of the stack be pushed to.
  - Suppose a list has no integer currently, then any of first integer from other lists is allowed to be popped and pushed to it.

#### III. Cost of different actions

- Cost is 1 for all actions,  $X_y \in A$  (uniform cost).

#### IV. Successor State (for each action)

- Suppose that an action,  $X_y$  is valid at a given state,  $s$  based on the constraints defined in (2) then the successor state,  $s'$  is defined as:
  - Let first integer of list X be  $k$ .
  - Then  $s'$  would be the state  $s$ , with  $k$  being popped from the front of list X, and pushed to the front of list Y.

## V. Objective

- Find the least cost path (a **sequence** of **operator** applied on **each intermediate states**) from the start state,  $s_0 = ([1,2...,N], [], [])$  to the goal state,  $s_g = ([], [], [1,2...,N])$ .

### Part 2

- For a heuristic to be admissible, the following must be satisfied:
$$h(n) \leq TC^*(n, g)$$
- Now, let  $h(n) = 2 * (\text{number of discs on tower C that is either smaller than discs on A or B}) + (\text{Total number of discs on tower A and B})$
- Justification of admissibility:
  - For any state, we have p discs on A, q discs on B, and r discs on C.
  - We also know the goal state is to have all discs being stacked appropriately (increasing size) on C.
  - Before moving any discs from A or B to C, we must first remove those discs with smaller size than those on A and B first (let this be m discs), this takes m actions.
  - Next, we perform relaxation on the original problem:
    - (1) Now we are allowed to remove any disc from any of the towers instead of only the top disc (but ONLY allowed to push to the top of the target tower) at any state.
    - (2) We have one extra tower, T and this tower is allowed hold any discs without considering the size.
  - With this relaxation, now we can remove the m discs on C to T, then always select the next largest discs from any of the A, B and T and place them in correct order on top C, this takes (m + p + q) actions.
  - Hence, total cost of the relaxed problem = m + (m + p + q) = 2 \* m + p + q
  - By the definition of original problem, this heuristic is trivially admissible as the total cost is lower bounded by the heuristic function.

### Part 3

- The priority of each state,  $n$  is based on the f-cost,  $f(n)$  which is defined as:

$$f(n) = g(n) + h(n)$$

That is, *sum of  $g(n)$ , cost from start state,  $s_0$  to current state,  $n$ , and  $h(n)$ , heuristic from current state,  $n$  to goal state,  $g$ .*

- Let the start state,  $S0_3$  be  $([1,2,3],[],[])$  and  $f(S0) = 0 + (2*0 + 3) = 3$
- Let the next valid successor states be:
- $S1_4^{S0}$  (after applying operator  $A_B$  on  $S0$ ):  $([2,3],[1],[])$  & f-cost =  $1 + (2*0 + 3) = 4$
- $S2_5^{S0}$  (after applying operator  $A_C$  on  $S0$ ):  $([2,3],[],[1])$  & f-cost =  $1 + (2*1 + 2) = 5$
- $S3_6^{S1}$  (after applying operator  $A_C$  on  $S1$ ):  $([3],[1],[2])$  & f-cost =  $2 + (2*1 + 2) = 6$
- $S0_5^{S1}$  (after applying operator  $B_A$  on  $S1$ ):  $([1,2,3],[],[])$  & f-cost =  $2 + (2*0 + 3) = 5$
- $S2_6^{S1}$  (after applying operator  $B_C$  on  $S1$ ):  $([2,3],[],[1])$  & f-cost =  $2 + (2*1 + 2) = 6$
- $S4_6^{S2}$  (after applying operator  $A_B$  on  $S2$ ):  $([3],[2],[1])$  & f-cost =  $2 + (2*1 + 2) = 6$
- $S0_5^{S2}$  (after applying operator  $C_A$  on  $S2$ ):  $([1,2,3],[],[])$  & f-cost =  $2 + (2*0 + 3) = 5$
- $S1_5^{S0}$  (after applying operator  $C_B$  on  $S2$ ):  $([2,3],[1],[])$  & f-cost =  $2 + (2*0 + 3) = 5$

...

First three steps of A\*star search in table representation:

Step #	Open List (Priority Queue)	Dequeue	Nodes to Add
1	$S0_3$	$S0_3$	$S1_4^{S0}, S2_5^{S0}$
2	$S1_4^{S0}, S2_5^{S0}$	$S1_4^{S0}$	$S3_6^{S1}, S0_5^{S1}, S2_6^{S1}$ * Do not re-open $S0$ in the close list * Do not replace $S2$ in the open list
3	$S2_5^{S0}, S3_6^{S1}$	$S2_5^{S0}$	$S4_6^{S2}, S0_5^{S2}, S1_5^{S0}$ * Do not re-open $S0$ in the close list * Do not re-open $S1$ in the close list
4	$S4_6^{S2}, S3_6^{S1}$	$S4_6^{S2}$	...

### QUESTION 3

#### Part (a)

The values for different states  $V^t$  at iteration t are given by:

<b>3</b>	1.564	1.710	1.842	<b>Food</b>
<b>2</b>	1.451	blocked	1.270	<b>Tiger</b>
<b>1</b>	1.327	1.231	1.162	0.844
	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>

Computation of values of different states  $V^{t+1}$  using Bellman Optimal Equation:

$$(1) Q^{t+1}(s, a) = \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V^t(s')]$$

$$(2) V^{t+1}(s) = \max_a Q^{t+1}(s, a)$$

The available actions at each state are = {N,S,E,W}, each representing  
N: North, S: South, E: East, W: West

#### For state (1,1):

$$Q^{t+1}((1,1), N)$$

$$= 0.7 * (-0.0025 + 0.95 * 1.451)$$

$$+ 0.15 * (-0.0025 + 0.95 * 1.327)$$

$$+ 0.15 * (-0.0025 + 0.95 * 1.231)$$

$$= 1.3269$$

$$Q^{t+1}((1,1), S)$$

$$= 0.7 * (-0.0025 + 0.95 * 1.327)$$

$$+ 0.15 * (-0.0025 + 0.95 * 1.327)$$

$$+ 0.15 * (-0.0025 + 0.95 * 1.231)$$

$$= 1.2445$$

$$Q^{t+1}((1,1), W)$$

$$= 0.7 * (-0.0025 + 0.95 * 1.327)$$

$$+ 0.15 * (-0.0025 + 0.95 * 1.451)$$

$$+ 0.15 * (-0.0025 + 0.95 * 1.327)$$

$$= 1.2758$$

$$Q^{t+1}((1,1), E)$$

$$= 0.7 * (-0.0025 + 0.95 * 1.231)$$

$$+ 0.15 * (-0.0025 + 0.95 * 1.451)$$

$$+ 0.15 * (-0.0025 + 0.95 * 1.327)$$

$$= 1.2120$$

$$V^{t+1}((1,1)) = 1.3269$$

#### For state (1,2):

$$Q^{t+1}((1,2), N)$$

$$= 0.7 * (-0.0025 + 0.95 * 1.231)$$

$$+ 0.15 * (-0.0025 + 0.95 * 1.327)$$

$$+ 0.15 * (-0.0025 + 0.95 * 1.162)$$

$$= 1.1708$$

$$\begin{aligned}
& Q^{t+1}((1,2), S) \\
&= 0.7 * (-0.0025 + 0.95 * 1.231) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.327) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.162) \\
&= 1.1708
\end{aligned}$$

$$\begin{aligned}
& Q^{t+1}((1,2), W) \\
&= 0.7 * (-0.0025 + 0.95 * 1.327) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.231) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.231) \\
&= 1.2308
\end{aligned}$$

$$\begin{aligned}
& Q^{t+1}((1,2), E) \\
&= 0.7 * (-0.0025 + 0.95 * 1.162) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.231) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.231) \\
&= 1.1211
\end{aligned}$$

$$V^{t+1}((1,2)) = 1.2308$$

**For state (1,3):**

$$\begin{aligned}
& Q^{t+1}((1,3), N) \\
&= 0.7 * (-0.0025 + 0.95 * 1.270) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.231) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 0.844) \\
&= 1.1377
\end{aligned}$$

$$\begin{aligned}
& Q^{t+1}((1,3), S) \\
&= 0.7 * (-0.0025 + 0.95 * 1.162) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.231) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 0.844) \\
&= 1.0659
\end{aligned}$$

$$\begin{aligned}
& Q^{t+1}((1,3), W) \\
&= 0.7 * (-0.0025 + 0.95 * 1.231) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.270) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.162) \\
&= 1.1627
\end{aligned}$$

$$\begin{aligned}
& Q^{t+1}((1,3), E) \\
&= 0.7 * (-0.0025 + 0.95 * 0.844) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.270) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.162) \\
&= 0.9053
\end{aligned}$$

$$V^{t+1}((1,3)) = 1.1627$$

**For state (1,4):**

$$\begin{aligned}
& Q^{t+1}((1,4), N) \\
&= 0.7 * (-1 + 0) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.162) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 0.844) \\
&= -0.4149
\end{aligned}$$

$$\begin{aligned}
& Q^{t+1}((1,4), S) \\
&= 0.7 * (-0.0025 + 0.95 * 0.844) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.162) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 0.844) \\
&= 0.8446
\end{aligned}$$

$$\begin{aligned}
& Q^{t+1}((1,4), W) \\
&= 0.7 * (-0.0025 + 0.95 * 1.162) \\
&\quad + 0.15 * (-1 + 0) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 0.844) \\
&= 0.7409
\end{aligned}$$

$$\begin{aligned}
& Q^{t+1}((1,4), E) \\
&= 0.7 * (-0.0025 + 0.95 * 0.844) \\
&\quad + 0.15 * (-1 + 0) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 0.844) \\
&= 0.5294
\end{aligned}$$

$$V^{t+1}((1,4)) = 0.8466$$

**For state (2,1):**

$$\begin{aligned}
& Q^{t+1}((2,1), N) \\
&= 0.7 * (-0.0025 + 0.95 * 1.564) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.451) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.451) \\
&= 1.4511
\end{aligned}$$

$$\begin{aligned}
& Q^{t+1}((2,1), S) \\
&= 0.7 * (-0.0025 + 0.95 * 1.327) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.451) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.451) \\
&= 1.2935
\end{aligned}$$

$$\begin{aligned}
& Q^{t+1}((2,1), W) \\
&= 0.7 * (-0.0025 + 0.95 * 1.451) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.564) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.327) \\
&= 1.3744
\end{aligned}$$

$$\begin{aligned}
& Q^{t+1}((2,1), E) \\
&= 0.7 * (-0.0025 + 0.95 * 1.451) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.564) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.327) \\
&= 1.3744
\end{aligned}$$

$$V^{t+1}((2,1)) = 1.4511$$

**For state (2,3):**

$$\begin{aligned}
& Q^{t+1}((2,3), N) \\
&= 0.7 * (-0.0025 + 0.95 * 1.842) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.270) \\
&\quad + 0.15 * (-1 + 0) \\
&= 1.2538
\end{aligned}$$



$$\begin{aligned}
& Q^{t+1}((2,3), S) \\
&= 0.7 * (-0.0025 + 0.95 * 1.162) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.270) \\
&\quad + 0.15 * (-1 + 0) \\
&= 0.8016
\end{aligned}$$

$$\begin{aligned}
& Q^{t+1}((2,3), W) \\
&= 0.7 * (-0.0025 + 0.95 * 1.270) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.842) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.162) \\
&= 1.2701
\end{aligned}$$

$$\begin{aligned}
& Q^{t+1}((2,3), E) \\
&= 0.7 * (-1 + 0) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.842) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.162) \\
&= -0.2727
\end{aligned}$$

$$V^{t+1}((2,3)) = 1.2701$$

**For state (3,1):**

$$\begin{aligned}
& Q^{t+1}((3,1), N) \\
&= 0.7 * (-0.0025 + 0.95 * 1.564) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.564) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.710) \\
&= 1.5041
\end{aligned}$$

$$\begin{aligned}
& Q^{t+1}((3,1), S) \\
&= 0.7 * (-0.0025 + 0.95 * 1.451) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.564) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.710) \\
&= 1.4290
\end{aligned}$$

$$\begin{aligned}
& Q^{t+1}((3,1), W) \\
&= 0.7 * (-0.0025 + 0.95 * 1.564) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.564) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.451) \\
&= 1.4672
\end{aligned}$$

$$\begin{aligned}
& Q^{t+1}((3,1), E) \\
&= 0.7 * (-0.0025 + 0.95 * 1.710) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.564) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.451) \\
&= 1.5643
\end{aligned}$$

$$V^{t+1}((3,1)) = 1.5643$$

**For state (3,2):**

$$\begin{aligned}
& Q^{t+1}((3,2), N) \\
&= 0.7 * (-0.0025 + 0.95 * 1.710) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.564) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.842) \\
&= 1.6200
\end{aligned}$$

$$\begin{aligned}
&Q^{t+1}((3,2), S) \\
&= 0.7 * (-0.0025 + 0.95 * 1.710) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.564) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.842) \\
&= 1.6200
\end{aligned}$$

$$\begin{aligned}
&Q^{t+1}((3,2), W) \\
&= 0.7 * (-0.0025 + 0.95 * 1.564) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.710) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.710) \\
&= 1.5249
\end{aligned}$$

$$\begin{aligned}
&Q^{t+1}((3,2), E) \\
&= 0.7 * (-0.0025 + 0.95 * 1.842) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.710) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.710) \\
&= 1.7098
\end{aligned}$$

$$V^{t+1}((3,2)) = 1.7098$$

**For state (3,3):**

$$\begin{aligned}
&Q^{t+1}((3,3), N) \\
&= 0.7 * (-0.0025 + 0.95 * 1.842) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.710) \\
&\quad + 0.15 * (2 + 0) \\
&= 1.7665
\end{aligned}$$

$$\begin{aligned}
&Q^{t+1}((3,3), S) \\
&= 0.7 * (-0.0025 + 0.95 * 1.270) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.710) \\
&\quad + 0.15 * (2 + 0) \\
&= 1.3861
\end{aligned}$$

$$\begin{aligned}
&Q^{t+1}((3,3), W) \\
&= 0.7 * (-0.0025 + 0.95 * 1.710) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.842) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.270) \\
&= 1.5781
\end{aligned}$$

$$\begin{aligned}
&Q^{t+1}((3,3), E) \\
&= 0.7 * (2 + 0) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.842) \\
&\quad + 0.15 * (-0.0025 + 0.95 * 1.270) \\
&= 1.8427
\end{aligned}$$

$$V^{t+1}((3,3)) = 1.8427$$

Hence, the numerical answer for each state of  $V^{t+1}$  is summarized as below:

<b>3</b>	1.5643	1.7098	1.8427	<b>Food</b>
<b>2</b>	1.4511	blocked	1.2701	<b>Tiger</b>
<b>1</b>	1.3269	1.2308	1.1627	0.8446
	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>

Part (b)

Optimal policy is chosen by :

$\pi^{t+1}(s) = \arg \max_{a \in \{N, S, W, S\}} Q^{t+1}(s, a)$ , which gives the best action (of that state) that maximize the expected utility over all states.

**For state (1,1):**

$$\pi^{t+1}((1,1)) = \arg \max_{a \in \{N, S, W, S\}} Q^{t+1}((1,1), a)$$

According to computation of  $Q^{t+1}((1,1), a)$ ,  $\forall a \in \{N, S, W, S\}$  in part (a), the action to be chosen for (1,1) should be N (North), corresponding to  $V^{t+1}((1,1)) = 1.3269$ .

**For state (1,2):**

$$\pi^{t+1}((1,2)) = \arg \max_{a \in \{N, S, W, S\}} Q^{t+1}((1,2), a)$$

According to computation of  $Q^{t+1}((1,2), a)$ ,  $\forall a \in \{N, S, W, S\}$  in part (a), the action to be chosen for (1,2) should be W (West), corresponding to  $V^{t+1}((1,2)) = 1.2308$ .

**For state (1,3):**

$$\pi^{t+1}((1,3)) = \arg \max_{a \in \{N, S, W, S\}} Q^{t+1}((1,3), a)$$

According to computation of  $Q^{t+1}((1,3), a)$ ,  $\forall a \in \{N, S, W, S\}$  in part (a), the action to be chosen for (1,3) should be W (West), corresponding to  $V^{t+1}((1,3)) = 1.1627$ .

**For state (1,4):**

$$\pi^{t+1}((1,4)) = \arg \max_{a \in \{N, S, W, S\}} Q^{t+1}((1,4), a)$$

According to computation of  $Q^{t+1}((1,4), a)$ ,  $\forall a \in \{N, S, W, S\}$  in part (a), the action to be chosen for (1,4) should be S (South), corresponding to  $V^{t+1}((1,4)) = 0.8466$ .

**For state (2,1):**

$$\pi^{t+1}((2,1)) = \arg \max_{a \in \{N, S, W, S\}} Q^{t+1}((2,1), a)$$

According to computation of  $Q^{t+1}((2,1), a)$ ,  $\forall a \in \{N, S, W, S\}$  in part (a), the action to be chosen for (2,1) should be N (North), corresponding to  $V^{t+1}((2,1)) = 1.4511$ .

**For state (2,3):**

$$\pi^{t+1}((2,3)) = \arg \max_{a \in \{N, S, W, S\}} Q^{t+1}((2,3), a)$$

According to computation of  $Q^{t+1}((2,3), a)$ ,  $\forall a \in \{N, S, W, S\}$  in part (a), the action to be chosen for (2,3) should be W (West), corresponding to  $V^{t+1}((2,3)) = 1.2701$ .

**For state (3,1):**

$$\pi^{t+1}((3,1)) = \arg \max_{a \in \{N, S, W, S\}} Q^{t+1}((3,1), a)$$

According to computation of  $Q^{t+1}((3,1), a)$ ,  $\forall a \in \{N, S, W, S\}$  in part (a), the action to be chosen for (3,1) should be E (East), corresponding to  $V^{t+1}((3,1)) = 1.5643$ .

**For state (3,2):**

$$\pi^{t+1}((3,2)) = \arg \max_{a \in \{N, S, W, S\}} Q^{t+1}((3,2), a)$$

According to computation of  $Q^{t+1}((3,2), a)$ ,  $\forall a \in \{N, S, W, S\}$  in part (a), the action to be chosen for (3,2) should be E (East), corresponding to  $V^{t+1}((3,2)) = 1.7098$ .

**For state (3,3):**

$$(9) \pi^{t+1}((3,3)) = \arg \max_{a \in \{N, S, W, S\}} Q^{t+1}((3,3), a)$$

According to computation of  $Q^{t+1}((3,3), a), \forall a \in \{N, S, W, E\}$  in part (a), the action to be chosen for (3,3) should be E(East), corresponding to  $V^{t+1}((3,3)) = 1.8427$ .

Hence, the final policy is summarized as follows:

<b>3</b>	East	East	East	<b>Food</b>
<b>2</b>	North	blocked	West	<b>Tiger</b>
<b>1</b>	North	West	West	South
	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>

## QUESTION 4

Part (a)

```
[8] # Import packages
import numpy as np
import gym
import random
import matplotlib.pyplot as plt
import seaborn as sn
%matplotlib inline

[10] # Create the Frozen Lake Environment
env = gym.make("FrozenLake8x8-v0")
env.render()

█ FFFFFFFF
FFFFFFF
FFTHFFFF
FFFFFFHF
FFFFFFF
FTHFFFF
FFTHHF
FFFFFFG

[11] # Determine state space
action_size = env.action_space.n
print("Action size: ", action_size)

# Determine action space
state_size = env.observation_space.n
print("State size: ", state_size)

# Create Q-table
qtable = np.zeros((state_size, action_size))
qtable_history = []
score_history = []
print("Table shape: ", qtable.shape)

Action size: 4
State size: 64
Table shape: (64, 4)

- Shows the size of the state space, action space
- Shows the appropriate data structure – Q-Table

[12] # Hyperparameters

gamma = 0.9          # Discount factor
total_episodes = 250000 # Total episodes
learning_rate = 0.8   # Learning rate
max_steps = 400       # Max steps per episode = action_size * state_size * 2

# Exploration parameters
epsilon = 1.0         # Exploration rate
max_epsilon = 1.0     # Exploration probability at start
min_epsilon = 0.001   # Minimum exploration probability
decay_rate = 0.00005  # Exponential decay rate for exploration prob
```

- Set the hyperparameters for Q-learning, and discount factor to 0.9

```
[13] # List of rewards
rewards = []

# Run the Q-Learning with pre-defined number of episodes
for episode in range(total_episodes):
    # for episode in range(10000):

        # Reset the environment in a new episode
        state = env.reset()

        step = 0

        done = False

        total_rewards = 0

        # In a episode (Start to Finish)
        for step in range(max_steps):

            # Choose an action, a in the current world state
            # Get a randomized number
            exp_exp_tradeoff = random.uniform(0, 1)

            # If this number > greater than epsilon : exploitation (taking the biggest Q value for this state)
            if exp_exp_tradeoff > epsilon:
                action = np.argmax(qtable[state,:])

            # Else doing a random choice : exploration
            else:
                action = env.action_space.sample()

            # Take the action (a) and observe the outcome state(s') and reward (r)
            new_state, reward, done, info = env.step(action)

            # Update Q(s,a) := Q(s,a) + lr [R(s,a) + gamma * max Q(s',a') - Q(s,a)]
            qtable[state, action] = qtable[state, action] + learning_rate * (reward + gamma * np.max(qtable[new_state, :]) - qtable[state, action])

            # Add the reward to total_rewards of this episode
            total_rewards += reward

            # Our new state is state
            state = new_state

            # If done (we're dead) : finish episode
            if done == True:
                break

        # Reduce epsilon after each iteration : exploitation over exploration
        epsilon = min_epsilon + (max_epsilon - min_epsilon)*np.exp(-decay_rate*episode)

    # Append the reward of this episode to reward list
    rewards.append(total_rewards)

    # Increase the episode_count
    episode_count = episode + 1
```

- Code that implements Q-learning
- **(For Part b)** Keep track of total accumulated reward for each episode (episode return), assigned to a list, *rewards*

## Part (b)

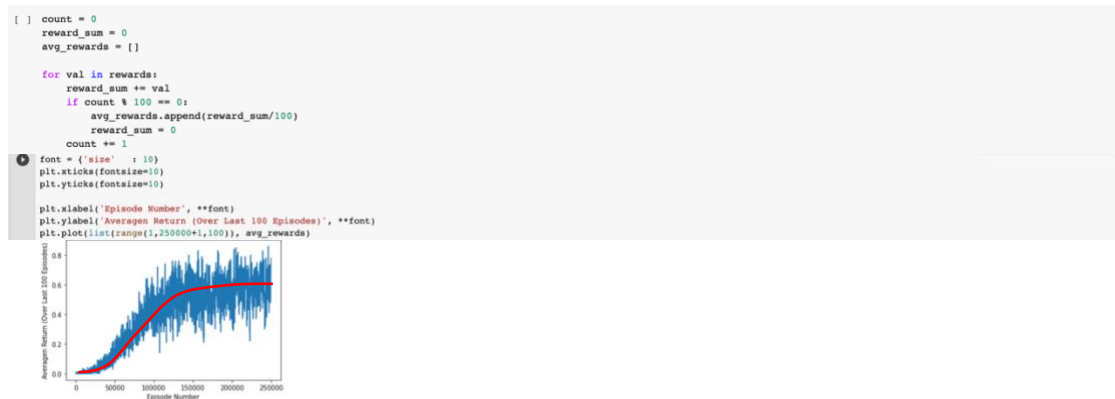
[illegible]

- Above shows the total accumulated reward for each episode

```
# Check the output
print("Reward over time: " + str(sum(rewards)/total_episodes))
print(qtable)
```

[illegible]

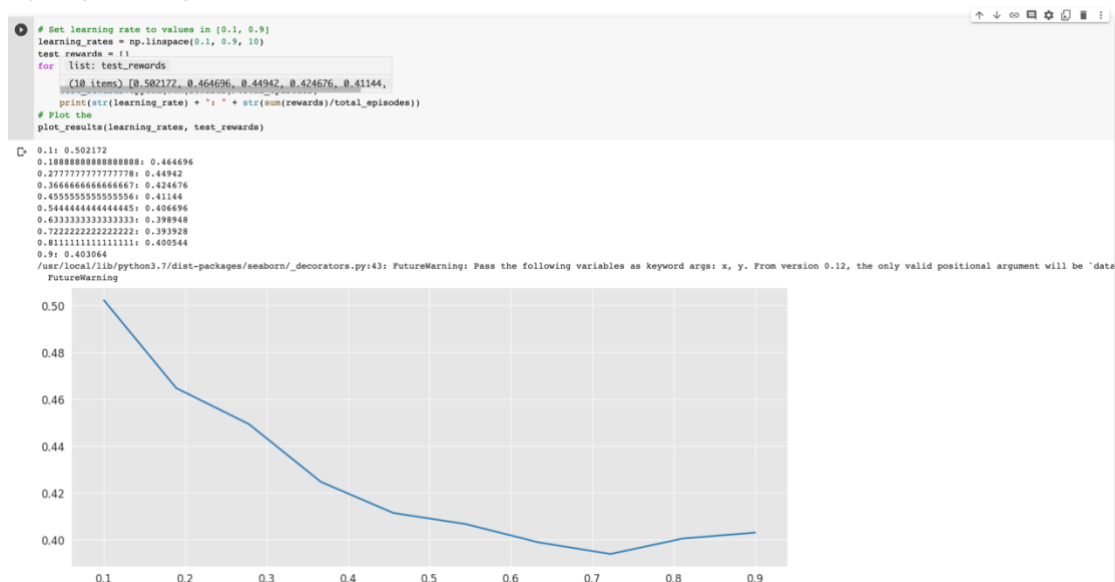
- Each of the values above is the Average Reward (Sum of Rewards of All Episodes / Total Number of Episodes)



- Plot the graph – average return (over the last 100 episodes while the agent is learning) on y-axis and episode number on x-axis
- The number of episodes is sufficient to ensure convergence as shown above

## Part (c)

Graph of Average Reward VS Learning Rate



- Plot the graph – Average Reward (Sum of Rewards of All Episodes / Total Number of Episodes) on y-axis and Learning Rate on x-axis
- As shown above, the Average Reward (as a metric to measure performance) drops when the Learning Rate increases
- According to [wikipedia](#), 'A too high learning rate will make the learning jump over minima but a too low learning rate will either take too long to converge or get stuck in an undesirable local minimum'
- In conclusion, after testing learning rates with 10 different values in [0.1,0.9], I believe 0.45556 is a decent learning rate to use, after considering both factors mentioned above
- As the average reward between using learning rate of 0.81111 and 0.45556 differ only approximately 0.01, and due to limitation of computation power of machine, I used **0.8** for this Q-Learning

### QUESTION 5

## Part (a)

```

import nltk

# import the stop words
nltk.download('stopwords')
from nltk.corpus import stopwords

# import the reuters dataset
nltk.download('reuters')
from nltk.corpus import reuters

file_loc = '/root/nltk_data/corpora/reuters.zip'
from zipfile import ZipFile
with ZipFile(file_loc, 'r') as z:
    z.extractall('/root/nltk_data/corpora/')

import numpy as np
from gensim.models import word2vec

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
[nltk_data] Downloading package reuters to /root/nltk_data...
[nltk_data] Package reuters is already up-to-date!

[ ] def load_data(category="money-fx"):

    """ Read files from the specified Reuter's category.
    Params:
        category (string): category name
    Return:
        list of lists, with words from each of the processed files/documents
    """

    files = reuters.fileids(category)

    list_words = []

    ### iterate over all documents/files from the reuters dataset
    for f in files:
        words = []

        ### iterate over all words of a document/file
        for w in list(reuters.words(f)):

            ## TODO: remove the stop words, convert a word (w) to the lowercase, and append to the list

            # (1) Convert w to lowercase and assign it to word
            word = w.lower()

            # (2) If the word is not a stopword, append it to list_words
            if word not in stopwords.words('english'):
                words.append(word)

        list_words.append(words)

    return list_words

# check a few samples of reuters corpus
reuters_corpus = load_data()
print(reuters_corpus[13])

[('bundesbank', 'allocates', '6', '.', '1', 'billion', 'marks', 'tender', 'bundesbank', 'accounted', 'bids', '6', '.', '1', 'billion', 'marks', 'today', 'tender', '28', '1', 'day', 'securities',

```

- Process words into the lowercase
- Perform stop words removal

## Part (b)

```

from nltk.corpus.reader.wordnet import lin_similarity
def distinct_words(corpus):
    """ get a list of distinct words for the corpus.
    Params:
        corpus (list of list of strings): corpus of documents
    Return:
        corpus_words (list of strings): list of distinct words across the corpus, sorted (using python 'sorted' function)
        num_corpus_words (integer): number of distinct words across the corpus
    """
    corpus_words = []
    num_corpus_words = -1

    # (1) Create a list of distinct words in the corpus, sorted using 'sorted' function
    for document in corpus:
        for word in document:
            corpus_words.append(word)

    corpus_words = list(dict.fromkeys(corpus_words))
    corpus_words = sorted(corpus_words)

    # (2) Find the number of distinct words across the corpus
    num_corpus_words = len(corpus_words)

    return corpus_words, num_corpus_words

words, num_words = distinct_words(reuters_corpus)
print('words (unique words): ', words)
print('num_words(unique words): ', num_words)

```

- Number of unique words (size of vocabulary) : 7157

```
[ ] def compute_co_occurrence_matrix(corpus, window_size=7):

    """ Compute co-occurrence matrix for the given corpus and window_size (default of 7).
    Params:
        corpus (list of list of strings): corpus of documents
        window_size (int): size of context window
    Return:
        M (numpy matrix of shape = [number of corpus words x number of corpus words]):
            Co-occurrence matrix of word counts.
            The ordering of the words in the rows/columns should be the same as the ordering of the words given by the distinct_words function.
        word2Ind (dict): dictionary that maps word to index (i.e. row/column number) for matrix M.
    """

    # (1) Create numpy matrix of shape = [number of corpus words x number of corpus words]
    M = np.array([[0 for i in range(num_words)] for i in range(num_words)])

    # (2) word2Ind (dict): dictionary that maps word to index (i.e. row/column number) for matrix M.
    word2Ind = {}
    idx = 0
    for str in words:
        word2Ind[str] = idx
        idx += 1

    # (3) Compute the co-occurrence matrix of word counts
    # Ordering of the words in the rows/columns should be the same as the ordering of the words given by the distinct_words function
    for document in corpus:
        for i in range(len(document)):
            for left_pointer in range(i-1, i-1-window_size, -1):
                if left_pointer < 0:
                    break
                M[word2Ind[document[i]]][word2Ind[document[left_pointer]]] += 1

            for right_pointer in range(i+1, i+1+window_size, 1):
                if right_pointer > (len(document) - 1):
                    break
                M[word2Ind[document[i]]][word2Ind[document[right_pointer]]] += 1

    return M, word2Ind

M, word2Ind = compute_co_occurrence_matrix(reuters_corpus)

[ ] #Check the dimension of Co-occurrence Matrix
print(len(M[0]))
print(len(M))
print(M)

7157
7157
[[240  0  1 ...  0  0  1]
 [ 0  0  0 ...  0  0  0]
 [ 1  0  0 ...  0  0  0]
 ...
 [ 0  0  0 ...  0  0  0]
 [ 0  0  0 ...  0  0  0]
 [ 1  0  0 ...  0  0  0]]
```

- Set the window size to 7 on either size of the word

## Part (c)

```
[ ] # -----
# Run SVD
# Note: This may take several minutes
# -----

# Perform SVD using
la = np.linalg

# Extract the first 75 components with highest variance
U, s, Vh = la.svd(M, full_matrices=False)

# Create an array of word embeddings (size of 75)
word_embedding_75dim = U[:, :75]

[ ] # Check dimension of word_embedding_75dim
print(len(word_embedding_75dim))
print(len(word_embedding_75dim[0]))

7157
75

[ ] # Check the entries in word_embedding_75dim
print(word_embedding_75dim[0])

[-0.09452258 -0.07720011 -0.12735456  0.05468409 -0.07876138 -0.04031832
 -0.22225215  0.09948283 -0.08399988  0.19380147 -0.31778021 -0.03067155
  0.02684725  0.13678115 -0.02323169 -0.03050785 -0.01286921 -0.06073952
  0.01030706 -0.07832958  0.04890193 -0.04293866 -0.03067227 -0.16071788
  0.11480831  0.11068655  0.16671438 -0.01143427 -0.39310589 -0.05267965
 -0.06456377 -0.02629117  0.10618979 -0.0766517  0.01791283 -0.25378902
 -0.02198051  0.01126623 -0.15405202 -0.1794979  0.11116424  0.28674949
  0.15448581 -0.14059331 -0.04498939  0.03704441 -0.10200549  0.09349053
 -0.00913624  0.03514359  0.00233887  0.03652442 -0.04197862 -0.10038659
  0.03754252  0.03747057 -0.05015903 -0.06717924  0.04048411  0.0154382
  0.02732327  0.04784751 -0.03990008 -0.06238556 -0.02026827  0.01941879
  0.04099045 -0.02403856  0.00823303  0.03884228  0.07178907 -0.03300502
  0.03370903  0.03295775 -0.02111137]
```

- Apply SVD and obtain word embedding of size 75



## Part (d)

```
from google.colab import drive
drive.mount('/content/drive')

# Creating the model and setting values for the various parameters
num_features = 75 # Set word vector dimensionality to 75 for consistency with previous implementation
min_word_count = 0 # Set minimum word count to 0 for consistency with previous implementation
num_workers = 4 # Number of parallel threads
context = 7 # Set context window size to 7 for consistency with previous implementation
downsampling = 1e-3 # (0.001) Downsample setting for frequent words

# Initializing the train model
from gensim.models import word2vec
print("Training model....")
model = word2vec.Word2Vec(reuters_corpus,\
                           workers=num_workers,\
                           size=num_features,\
                           min_count=min_word_count,\
                           window=context,\
                           sample=downsampling)

# To make the model memory efficient
model.init_sims(replace=True)

# Saving the model for later use ( Word2Vec.load() )
model_name = "75_features_Doinwords_7context"
model_path = f"/content/drive/MyDrive/colab_notebooks/asn2/{model_name}"
model.save(model_path)
print("Model saved")

Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount('/content/drive', force_remount=True).
Training model....
model saved
```

- Generate word embedding of size 75 using Word2Vec.pybn on the same Reuters dataset

```
from sklearn.metrics.pairwise import cosine_similarity

def svd_most_similar(query_word, n=10):
    """ return 'n' most similar words of a query word using the SVD word embeddings similar to word2vec's most_similar
    Params:
        query_word (strings): a query word
    Returns:
        most_similar (list of strings): the list of 'n' most similar words
    """
    cosine_sim = cosine_similarity(word_embedding_75dim, np.array([word_embedding_75dim[word2ind(query_word)]]))
    cosine_sim = cosine_sim.flatten()
    indices = (-cosine_sim).argsort()[1:n+1]
    most_similar = []
    for i in indices:
        most_similar.append((words[i], cosine_sim[i]))
    for val in most_similar:
        print(val)
```

- To compare and determine which method (Word2Vec OR SVD) works better, a few examples are shown below:

```
[95] model = word2vec.Word2Vec.load(model_path)
model.vv.most_similar('money')
```

```
[('assistance', 0.9967169761657715),
 ('operated', 0.996709910256958),
 ('revised', 0.9957497715950012),
 ('forecast', 0.9947110414505005),
 ('provided', 0.9942721128463745),
 ('700', 0.9939642548561096),
 ('350', 0.993422269821167),
 ('receives', 0.9923985930938721),
 ('compares', 0.9910964369733865),
 ('market', 0.9906169772148132)]
```

```
[98] svd_most_similar('money')
```

```
[('gets', 0.4576669125840092),
 ('442', 0.4521483187068887),
 ('given', 0.4492420270849863),
 ('210', 0.43187222984079865),
 ('205', 0.4255437853248652),
 ('receives', 0.4222603958836185),
 ('186', 0.39136546939813225),
 ('injects', 0.3900399422575899),
 ('570', 0.3863047941647073),
 ('166', 0.3796084071326851)]
```

```
[101] model.vv.most_similar('bank')
```

```
[('help', 0.9957151412963867),
 ('given', 0.9939113259315491),
 ('market', 0.9924650192260742),
 ('affecting', 0.9912904500961304),
 ('afternoon', 0.9896146655082703),
 ('700', 0.986505315856934),
 ('invited', 0.9858958721160889),
 ('taken', 0.9856358766555786),
 ('191', 0.9854189157485962),
 ('750', 0.9851484298706055)]
```

```
[102] svd_most_similar('bank')
```

```
[('stereotype', 0.3016644191084184),
 ('tide', 0.2533900827092676),
 ('jokes', 0.2401320327752718),
 ('sportsman', 0.23768650897948743),
 ('beez', 0.233881858995177),
 ('glass', 0.23212708691109467),
 ('scored', 0.23198434157524173),
 ('sufficed', 0.22786686713750873),
 ('megarditch', 0.22402487622726294),
 ('like', 0.22377866923413883)]
```

- Based on the observation above, on overall, Word2Vec approach works better for this reuters\_corpus:
  - o For the word 'money', 5 out of 10 of the returned results of SVD approach are numerical values, which is relatively less meaningful for this case as compared to Word2Vec approach.

- Another possible reason is the corpus contains relative high number of non-English words. Although these special characters may be meaningful to certain extend, but with word counts approaching 10% of the entire corpus may affects, and possibly reduce the accuracy of SVD computation approach.

- Finish -