

Introduction To Reinforcement Learning

Project 2

Koray Emre Şenel - 150117037

Problem:

In this project, Q Learning algorithm is implemented.

Q-Learning:

For Q-learning this algorithm is used and shown which part of it, is done in the code with comments.

Algorithm for Q-learning

- *Initialize $Q(s,a)$ arbitrarily*
 - *Repeat for each episode*
 - *Initialize s ;*
 - *Repeat for each step of episode*
 - *Choose a from s using policy derived from Q (e.g., ϵ -greedy)*
 - *Take action a , observe reward r , and next state s'*
- $$Q(s, a) = Q(s, a) + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$
- *$s \leftarrow s'$;*
 - *Until s is terminal*
- *Until convergence occurs*

This is the main function. In this function initialization of start state and repeat for each episode part is done.

```
public static void main(String[] args) {  
    int i = 0;  
    int counter ;  
    int iteration = 100000000;  
    int startState = 0;  
    for (int i = 0; i < iteration; i++) {  
        counter = episodes(startState, seeRoad: false);  
    }  
}
```

In the episodes function, we have a loop with a flag condition where this flag is true until load, unload and returning to start state is done. The corresponding parts are commented in code.

```
208 static int episodes( int startState, boolean seeRoad) {  
209     int counter = 0;  
210     int agentState = startState; //Initialize s  
211     boolean flag = true;  
212     isDone = false;  
213     isLoading = false;  
214     boolean isLoadingSwap = true;  
215     boolean isDoneSwap = true;  
216     while (flag) { //Repeat for each state of episode  
217         counter++;  
218         int nextState;  
219         int act ;  
220         act = findAct(agentState); //Choose a from s using policy derived from epsilon greedy.  
221         nextState = takeAction(act, agentState); //Take action and observe reward r and next state s'  
222         if(nextState != agentState){  
223             lastAct = act;  
224             if(queue.size() > 19){  
225                 queue.remove();  
226             }  
227             queue.add(agentState);  
228         }  
    }  
}
```

The epsilon greedy method, a random number is generated between (0,1), if it is smaller than epsilon the random action is selected, else action with max Q value is selected.

```
static int findAct(int currentState){
    double rand = random.nextDouble();
    int act ;
    if (epsilon > rand) {
        if(currentState == loadState || currentState == unloadState){
            act = random.nextInt( bound: 5);
        }
        else act = random.nextInt( bound: 4);
    }
    else {
        act = actQ(currentState);
    }
    return act;
}
```

This part of code makes the Q table update for s,a.

```
double q = Q[agentState][act]; //s states Q value
double max = maxQ(nextState); //Max actions Q value for s'
double value = q + alpha * (reward + gamma * max - q); //New value of Q(s,a)
setQ(value, agentState, act); //Update function
```

Next state set as agent state.

```
agentState = nextState;
```

Check for terminal state and reward.

```
if(isDone && nextState == startState){  
    reward = 1000;  
    flag = false;  
}
```

In converge part, after 1000 episodes each 20 episode is taken and if variance of it is smaller than 0.1 execution terminated the result is printed.

```
static boolean isLearned(int i, int counter) {  
    int[] arr ;  
    arr = stepQueue.stream().mapToInt(Integer::intValue).toArray();  
    double sumOfSample = 0;  
    double variance;  
    double sampleMean;  
    for (int j = 0; j < i; j++)  
        sumOfSample += arr[j];  
    sampleMean = sumOfSample / i;  
    double sumVari = 0;  
    for (int j = 0; j < i; j++) {  
        sumVari += Math.pow(((double) arr[j] - sampleMean), 2);  
    }  
    variance = sumVari / (double) (i - 1);  
    if (variance < 0.1) {  
        System.out.println("solved at " + counter + " steps");  
        return true;  
    } else return false;  
}
```

What else used to solve this problem:

For encouraging agent to make load and unload operations, this operations gives 0 reward instead of -0.1. Also to check the correct road usage to load, looked the last 20 moves of the agent and check 3 of them have the correct road. If this road is not used, load is not working. For cross walls, the last successful action is stored, and used to determine which way it came to cross walls, action is done with that information.

Results of Q-learning:

With the settings of $\alpha = 0.1$, $\gamma = 0.9$, $\epsilon = 0.3$ (decreasing each episode) it converges to 46 steps (expected result) in approximately 6000 episodes.

```

Main x
↑ train 5921 solved at 46 steps
↓ train 5922 solved at 46 steps
↺ train 5923 solved at 46 steps
↻ train 5924 solved at 46 steps
⇅ train 5925 solved at 46 steps
⇅ train 5926 solved at 46 steps
⇅ train 5927 solved at 46 steps
⇅ train 5928 solved at 46 steps
⇅ train 5929 solved at 46 steps
⇅ train 5930 solved at 46 steps
⇅ train 5931 solved at 46 steps
⇅ train 5932 solved at 46 steps
⇅ train 5933 solved at 46 steps
solved at 46 steps
0->10->20->30->40->41->51->61->62->72->73->83->84->85->86->87->88->98->99->
99->89->79->69->59->49->39->29->19->9->9->19->29->28->38->48->47->46->45->44->34->24->14->13->12->2->1->End
Process finished with exit code 0
```

With the settings of $\alpha = 0.1$, $\gamma = 0.8$, $\epsilon = 0.4$ (decreasing each episode) it converges to 46 steps (expected result) in approximately 8000 episodes.

```

Main x
↑ train 7952 solved at 46 steps
↓ train 7953 solved at 46 steps
↺ train 7954 solved at 46 steps
↻ train 7955 solved at 46 steps
⇅ train 7956 solved at 46 steps
⇅ train 7957 solved at 46 steps
solved at 46 steps
0->10->20->30->40->41->51->52->62->63->73->74->84->85->95->96->97->98->99->
99->89->79->69->59->49->39->29->19->9->9->19->18->28->38->48->47->46->45->44->43->33->23->13->12->11->10->End
Process finished with exit code 0
```

With the settings of $\alpha = 0.9$, $\gamma = 0.3$, $\epsilon = 0.6$ (decreasing each episode) it converges to 46 steps (expected result) in approximately 12000 episodes.

```

Main x
↑ train 11846 solved at 46 steps
↓ train 11847 solved at 46 steps
↺ train 11848 solved at 46 steps
↻ train 11849 solved at 46 steps
⇓ train 11850 solved at 46 steps
⇓ train 11851 solved at 46 steps
⇓ solved at 46 steps
0->10->11->12->11->21->31->41->51->52->62->72->73->74->84->94->95->96->97->98->99->
99->89->79->69->59->49->39->29->19->9->9->8->18->28->38->48->47->46->45->44->43->33->32->22->12->2->1->End
Process finished with exit code 0

```

With the settings of $\alpha = 0.1$, $\gamma = 0.9$, $\epsilon = 0.1$ (decreasing each episode) it converges to 46 steps (expected result) in approximately 2000 episodes. This gives the best result in my experiments.

```

↑ train 1850 solved at 46 steps
↓ train 1851 solved at 46 steps
↺ train 1852 solved at 46 steps
↻ train 1853 solved at 46 steps
⇓ train 1854 solved at 46 steps
⇓ train 1855 solved at 46 steps
⇓ train 1856 solved at 46 steps
⇓ solved at 46 steps
0->10->20->30->31->41->51->61->62->72->73->83->84->85->86->87->97->98->98->99->
99->89->79->69->59->49->39->29->19->9->9->19->18->28->38->48->47->46->45->44->34->33->32->22->21->11->1->End
Process finished with exit code 0

```