# Problem 1

## Question1.

In this question use vector-form Policy Iteration technique to find the optimal policy and the optimal state values. Choose action Left in all states as your initial policy for the Policy Iteration method and all-zero initial state values for the Value Iteration method.

In the Base Scenario factor, because the probability of selecting the optimal action is high , so the algorithm is likely to converge more rapidly to the optimal policy. This is because the agent is more likely to select the action with the highest expected reward, which leads to a more efficient search for the optimal policy. In the Large Stochasticity Scenario factor, because the optima policy would be selected by a small probability, which will result in a suboptimal policies. So compare to the Large Stochasticity Scenario, it is easy for Base Scenario factor to get the optimal policy.
In the Small Discount Factor Scenario , because it has the lower gamma compare to the other two factors, so the agent will give higher priority to immediate rewards over long-term rewards, in this situation the agent would rather take risky or uncertain actions, so in the give gamma of the question, the agent was trapped in the loop.

The plots of Policy Iteration values,actions and path are as follow:
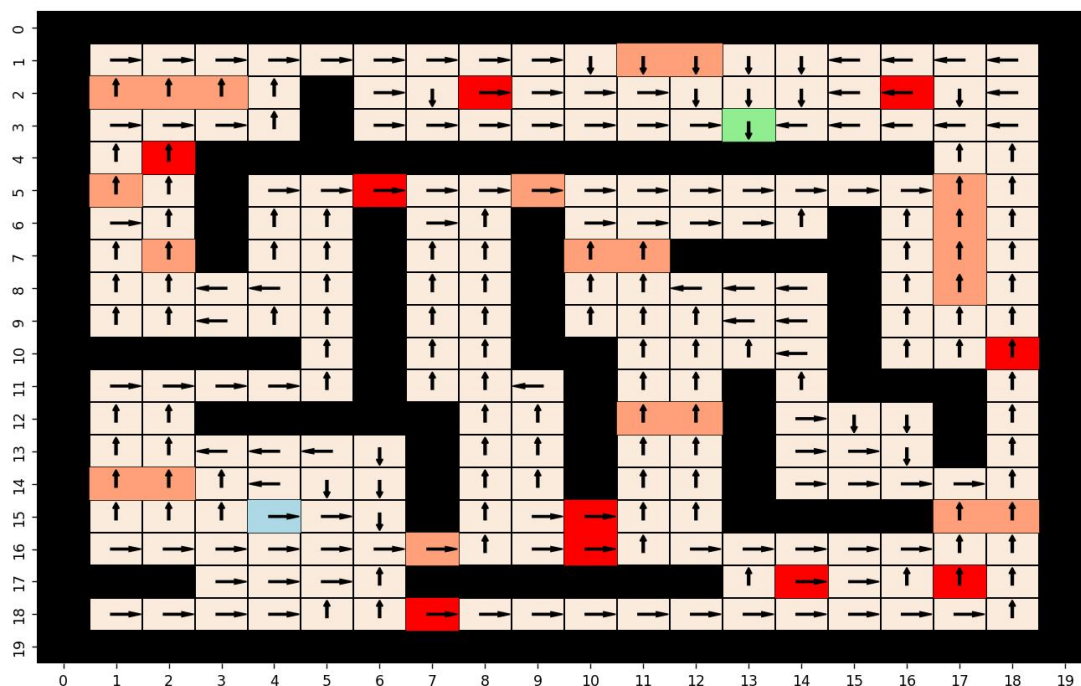
a. P =0.02，γ =0.95，θ =0.01



Figure 1.1 The optimal policy for each state for the base scenario under the Policy Iteration method.
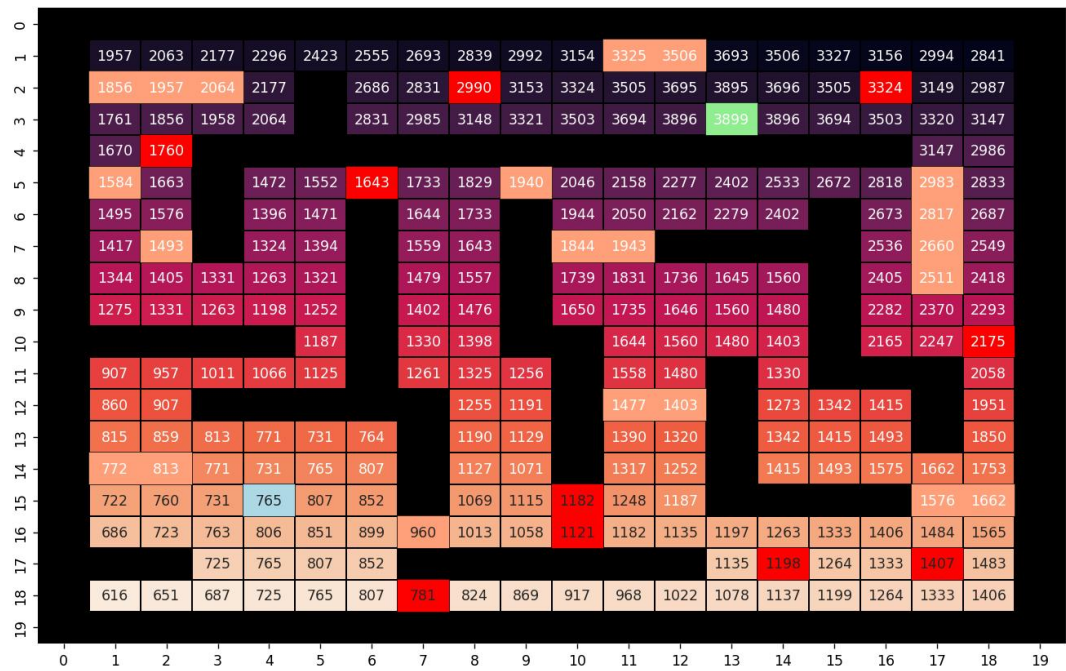
**Figure 1.2 The optimal value function values for each state for the base scenario under the Policy Iteration method**
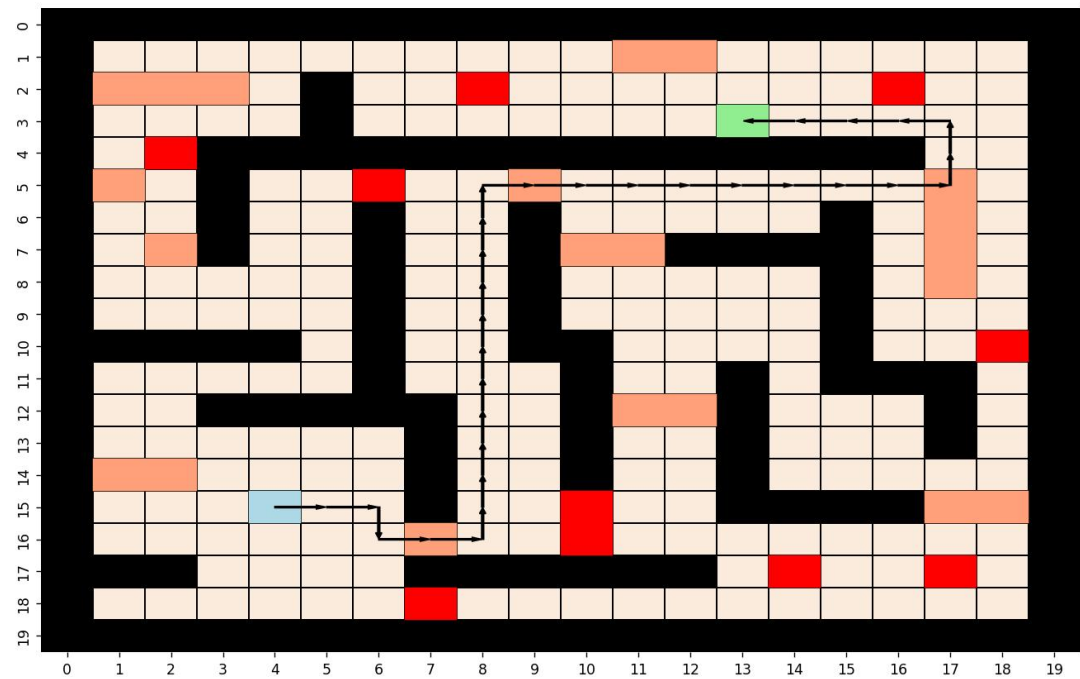


**Figure 1.3 The optimal path from the start state to the end state for the base scenario under the Policy Iteration method.**
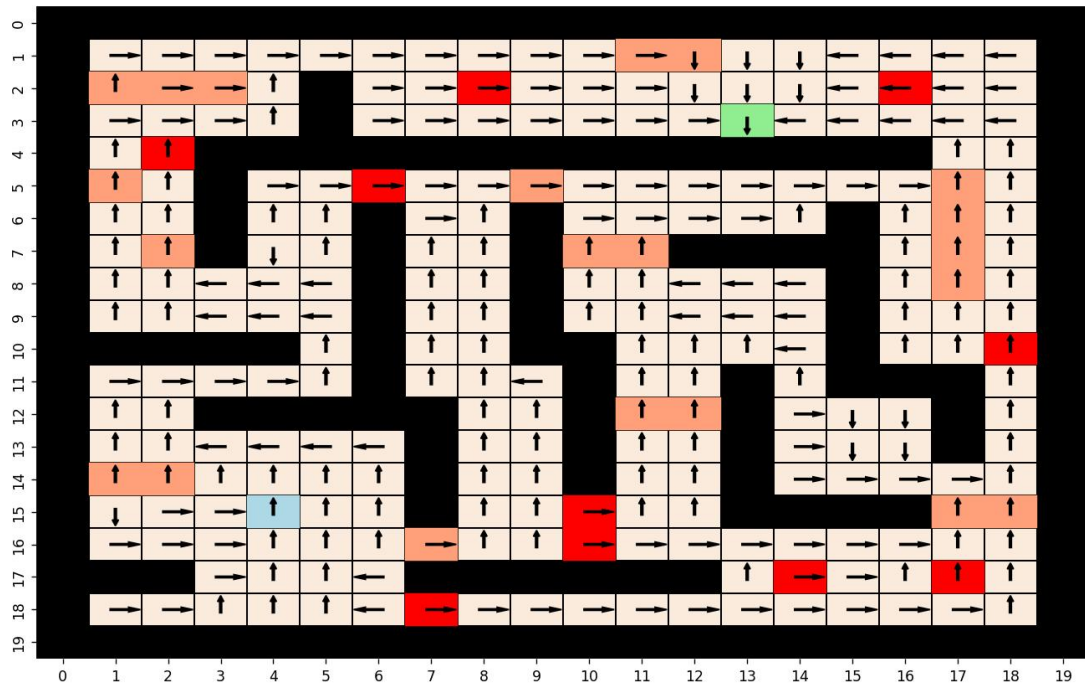
**Figure 1.4The optimal policy for each state for the base scenario under the Policy Iteration method.**
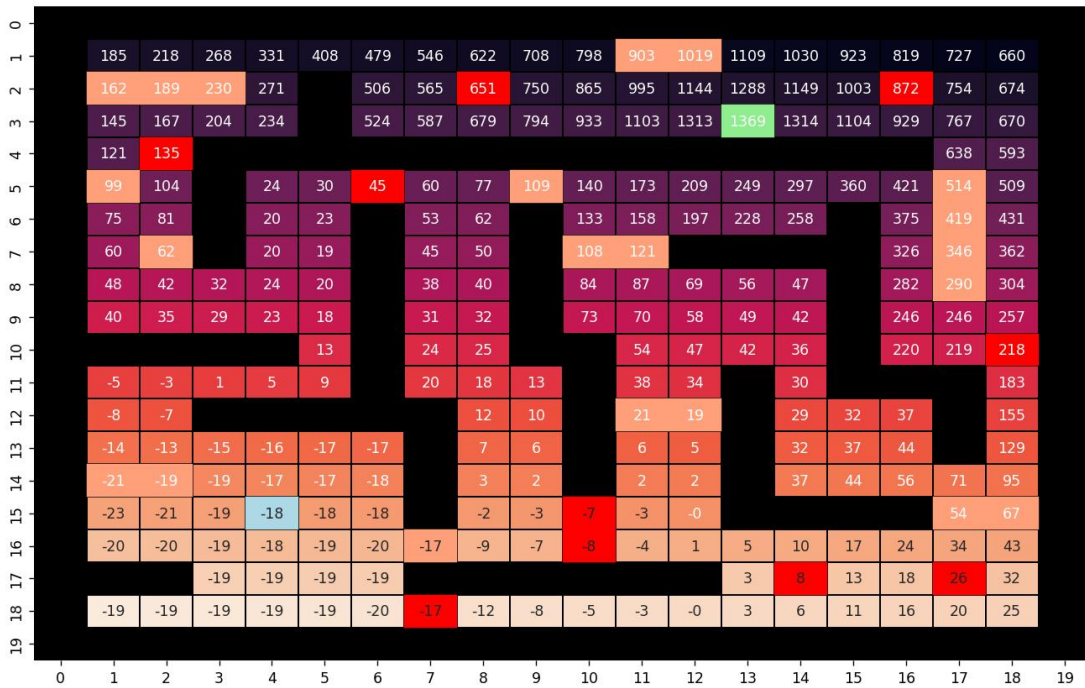


**Figure 1.5The optimal value function values for each state for the base scenario under the Policy Iteration method**

**Figure 1.6The optimal path from the start state to the end state for the base scenario under the Policy Iteration method.**

c. P =0.02, γ =0.55，θ =0.01



**Figure 1.7 The optimal policy for each state for the base scenario under the Policy Iteration method.**

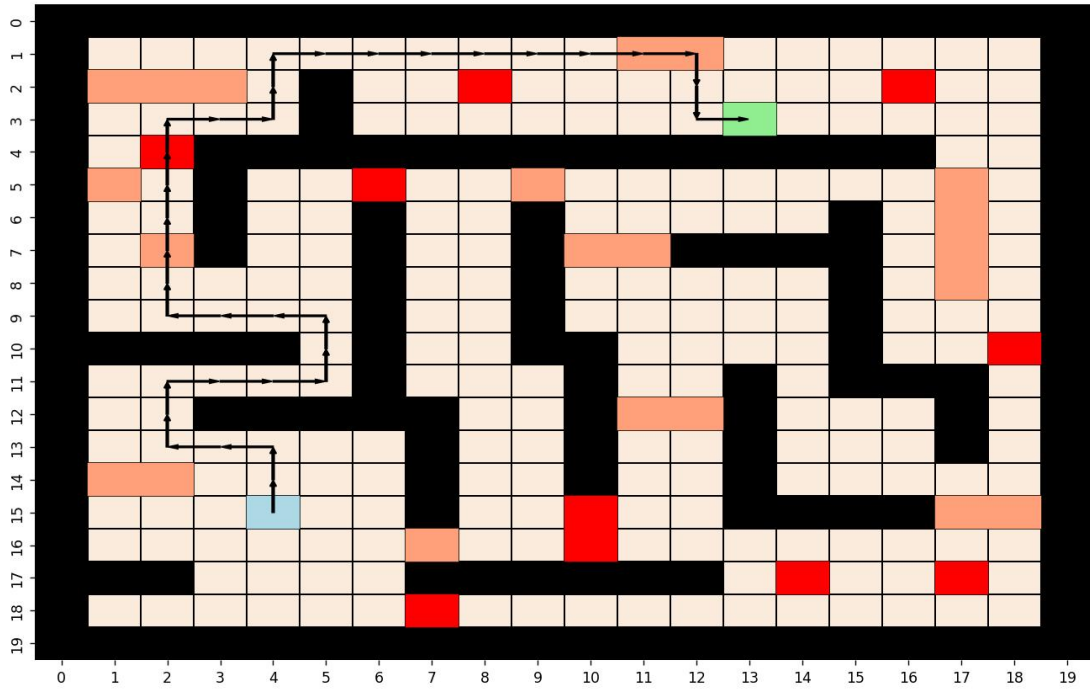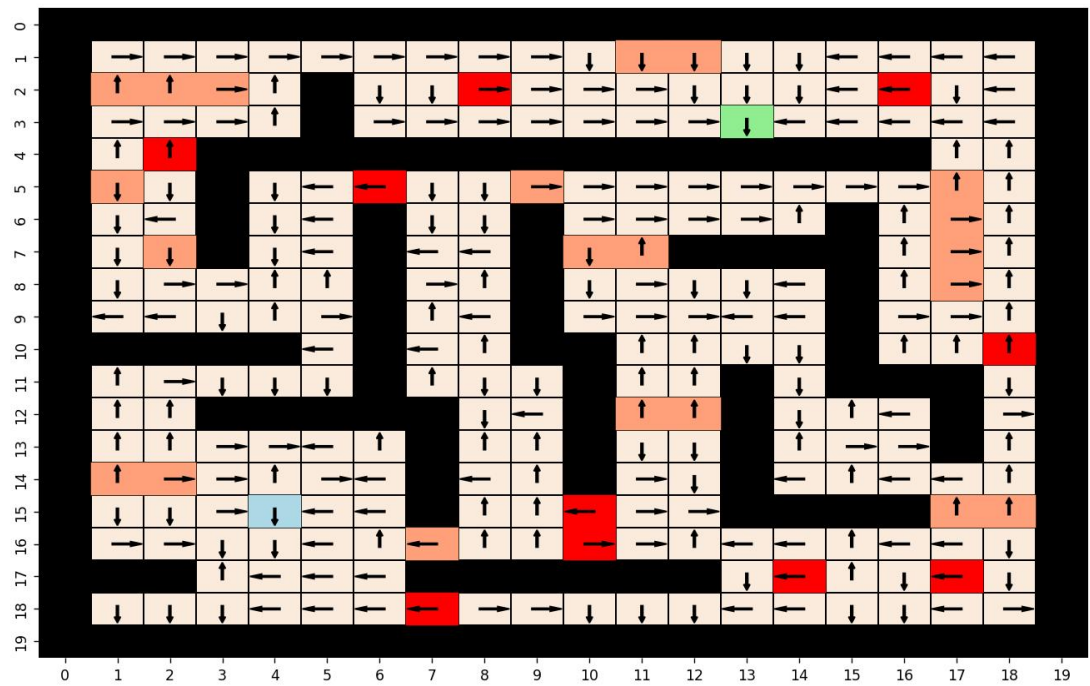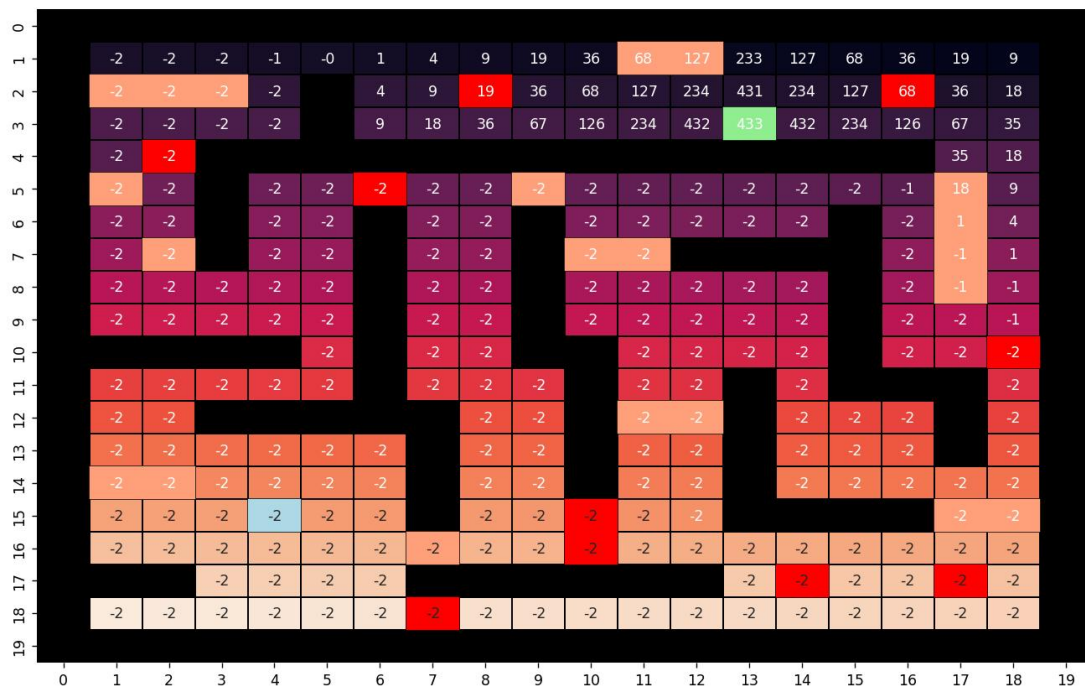**Figure 1.8The optimal value function values for each state for the base scenario under the Policy Iteration method**
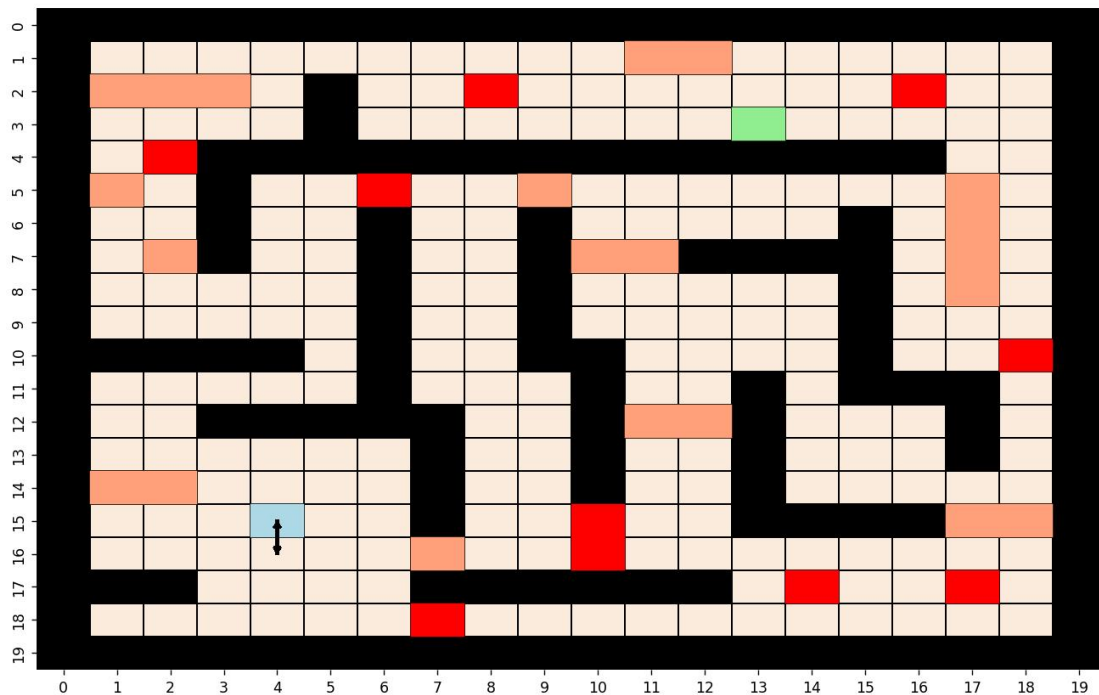


**.Figure 1.9The optimal path from the start state to the end state for the base scenario under the Policy Iteration method.**

## Question2.

Repeat the same process as part 1 with the vector-form Value Iteration method. Compare Value Iteration and Policy Iteration techniques for this problem

According to the different factors outcome, using Base Scenario factor and Large Stochasiticity Scenario factor in Value Iteration would take more iterations than Policy Iteration. Which because that Policy Iteration is better suited for situations where the optimal policy is needed quickly. So contract to the Policy Iteration, if the optimal policy is not required immediately, and MDP also has the large state space, value iteration would be better to use.

The plots of Policy Iteration values,actions and path are as follow:

a. P =0.02，γ =0.95， θ =0.01



Figure 2.1 The optimal policy for each state for the base scenario under the Value Iteration method



Figure 2.2 The optimal value function values for each state for the base scenario under the Value Iteration method

**Figure 2.3 The optimal path from the start state to the end state for the base scenario under the Value Iteration method.**
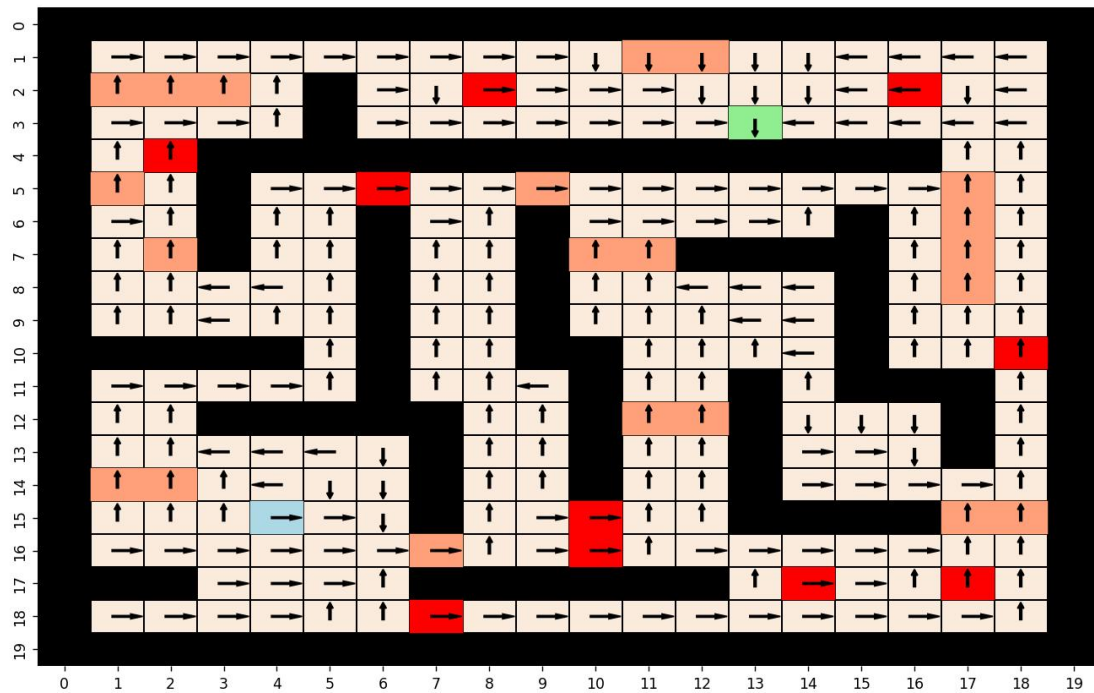
b. P =0.5，γ =0.95， θ =0.01



**Figure 2.4 The optimal policy for each state for the base scenario under the Value Iteration method**
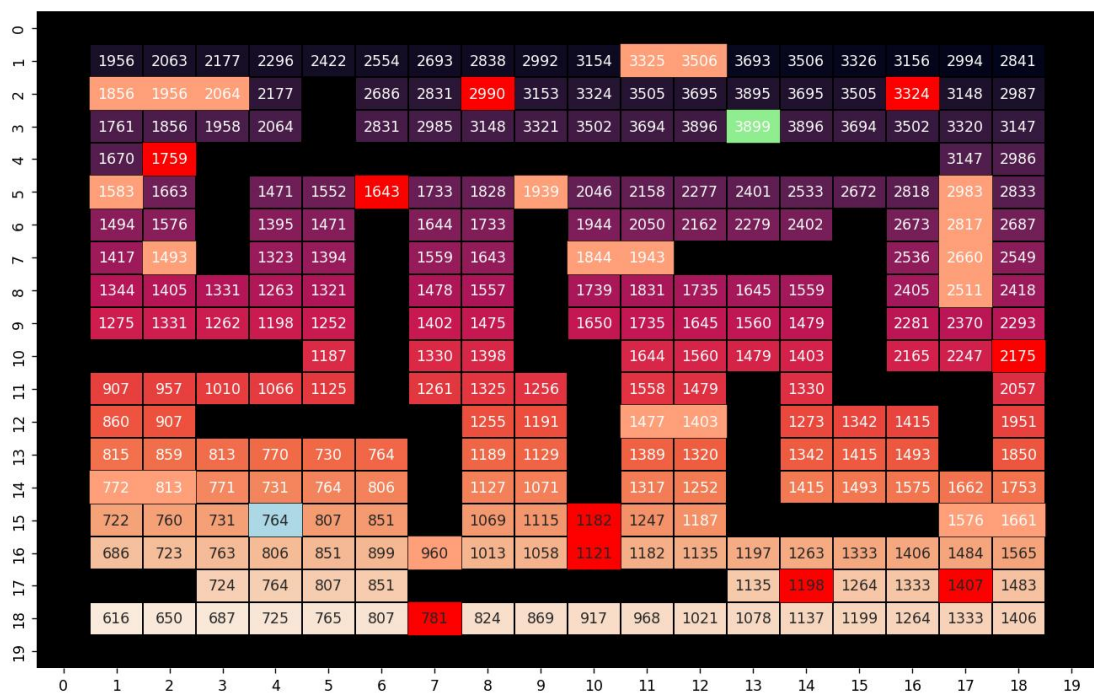
**Figure 2.5 The optimal value function values for each state for the base scenario under the Value Iteration method**



**Figure 2.6 The optimal path from the start state to the end state for the base scenario under the Value Iteration method.**

**Figure 2.7 The optimal policy for each state for the base scenario under the Value Iteration method**



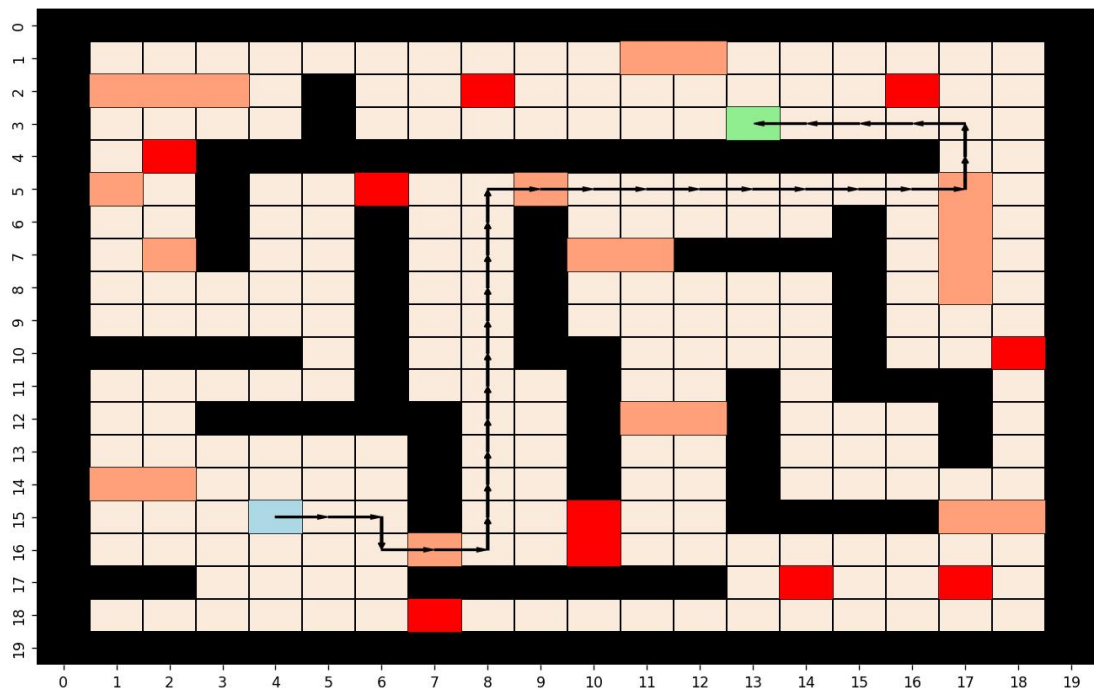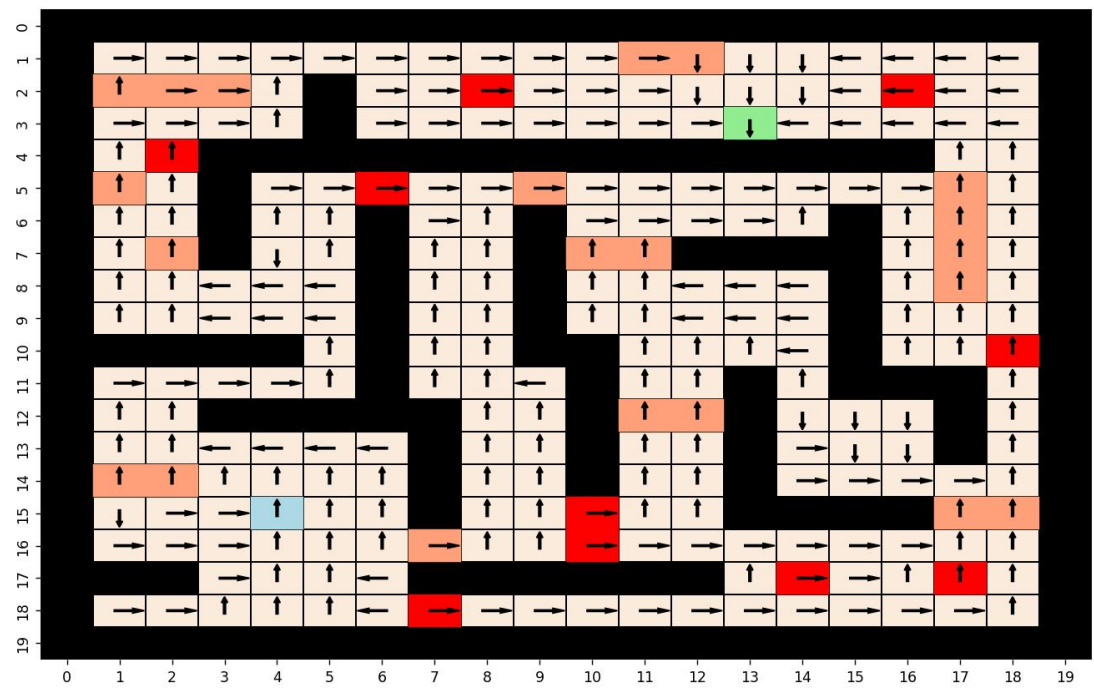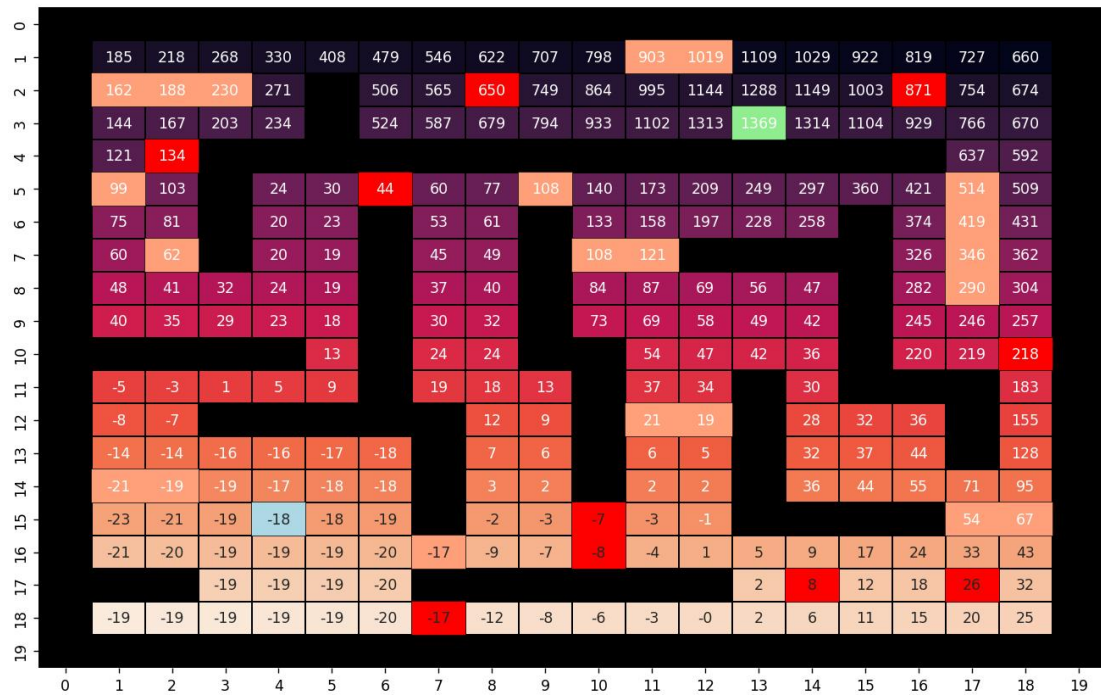**Figure 2.8 The optimal value function values for each state for the base scenario under the Value Iteration method**
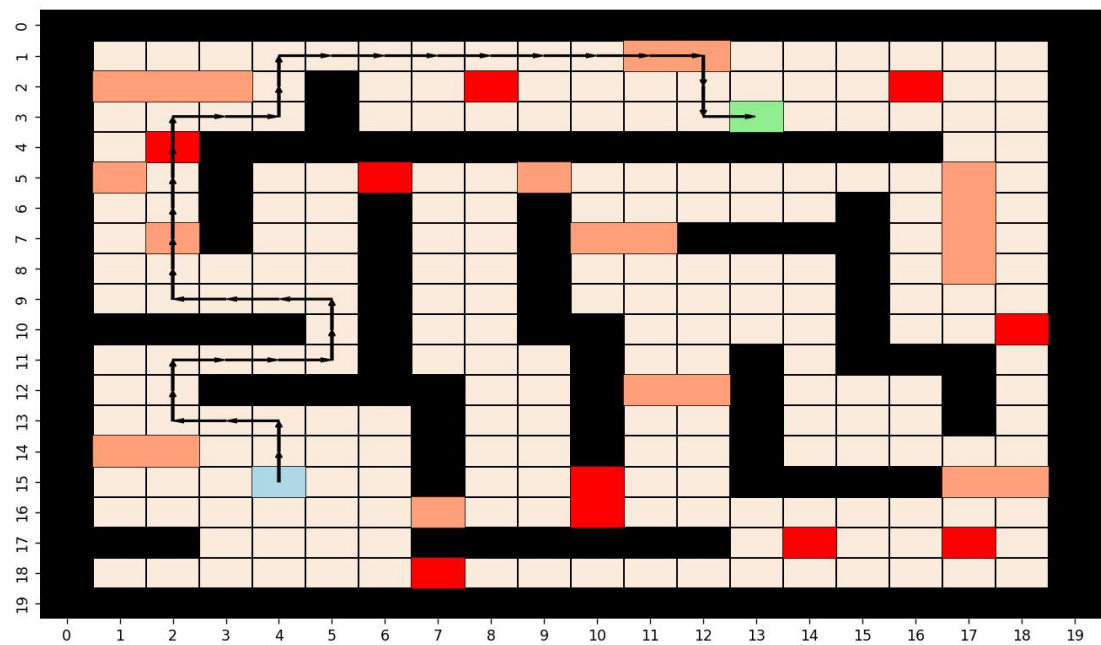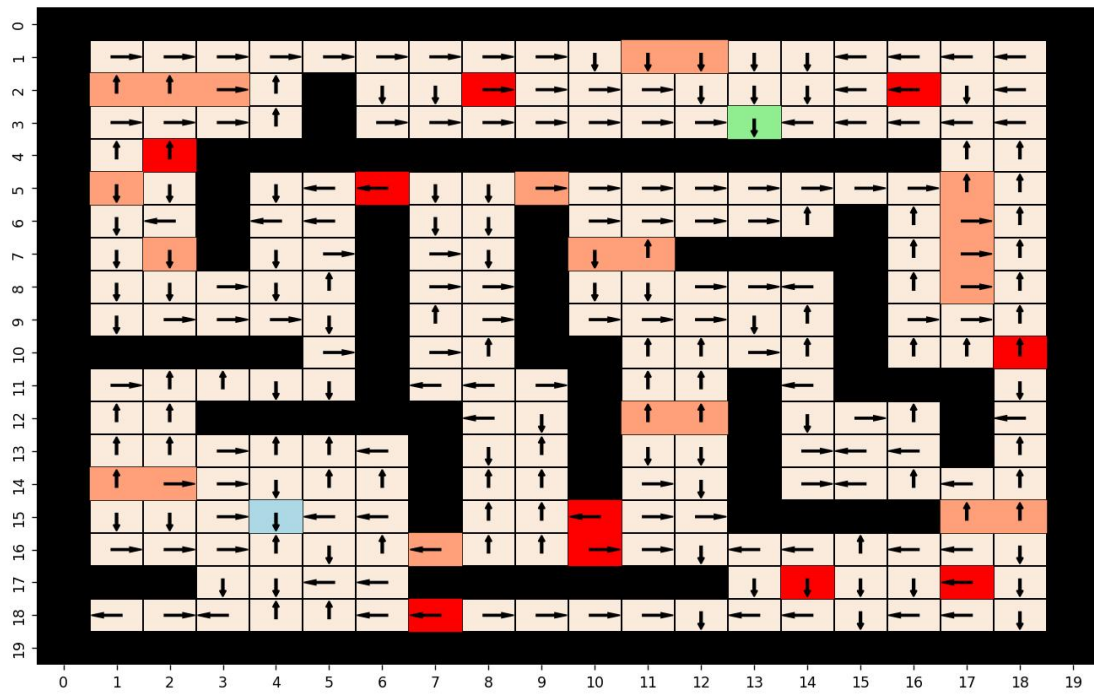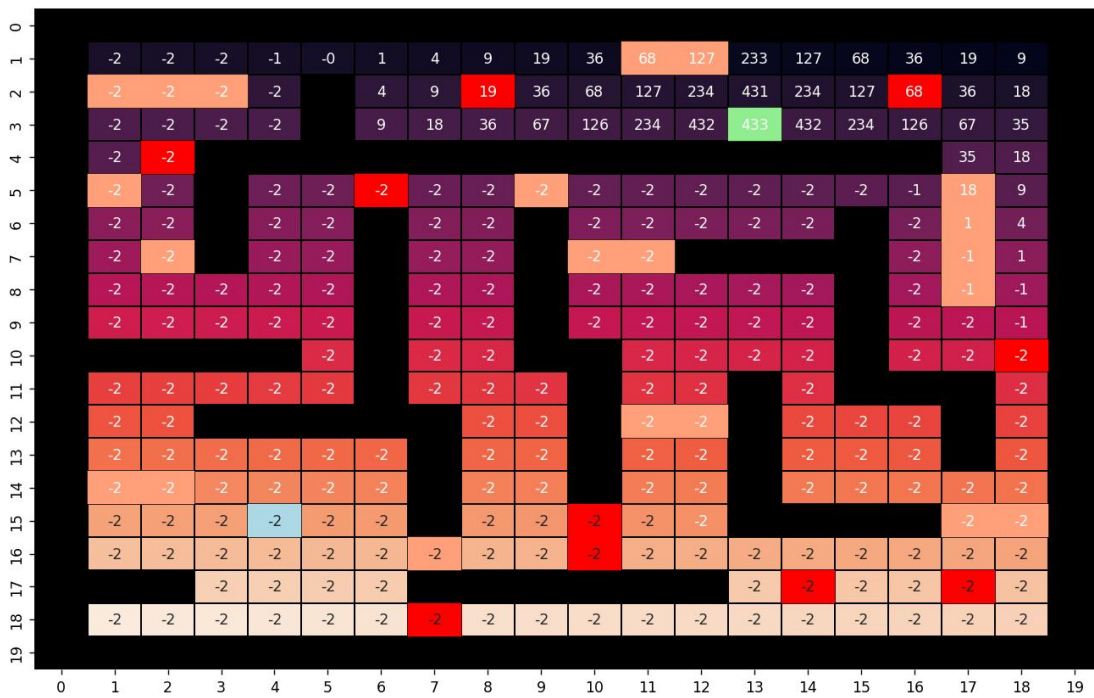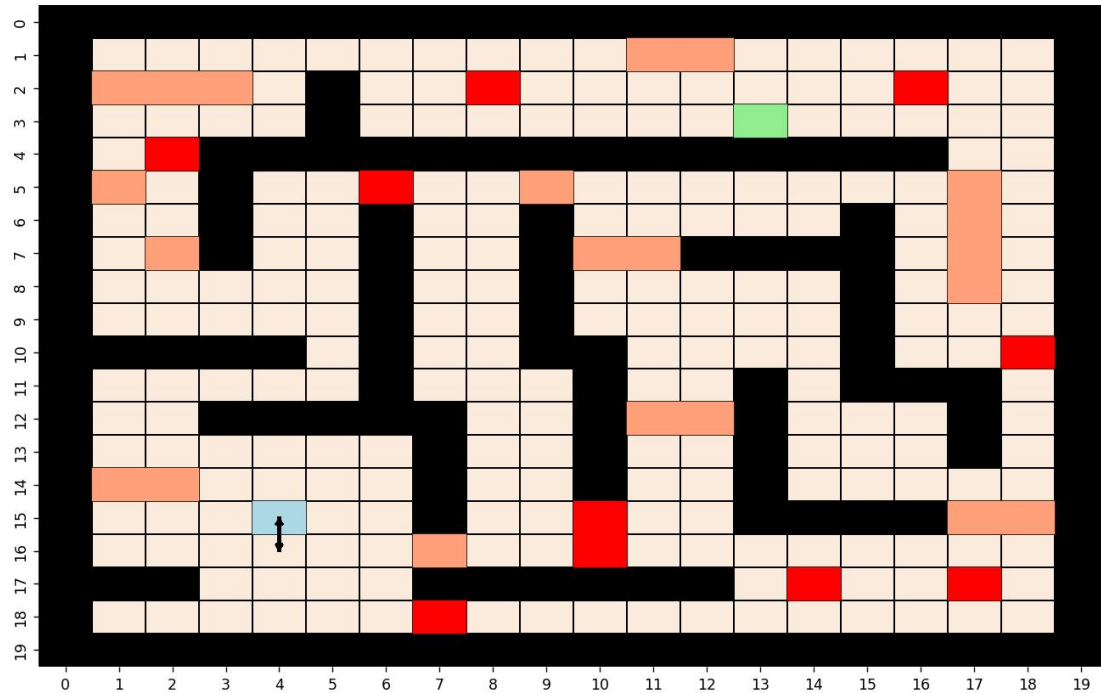
**Figure 2.9 The optimal path from the start state to the end state for the base scenario under the Value Iteration method.**

# Problem 2

## Part A

For p=0.05 use matrix-form Value Iteration method for computing the optimal control policy $\pi^*$, which is a vector of size 16, specifying the best action at any given state. Use all-zero initial state values and $\theta$ =0.01. Compare the AvgA under the obtained optimal control policy with AvgA under no control policy(taking a-in all states)

.

p=0.05

```
values= [[255.45551755 255.45551755 259.7727771  259.7727771  264.62085953
  264.62085953 264.62085953 264.62085953 260.13047744 255.45551755
  255.45551755 255.45551755 264.62085953 260.13047744 264.62085953
  264.62085953]]
policy= [2, 2, 2, 2, 2, 2, 2, 2, 3, 2, 2, 2, 4, 2, 2, 2]
```

```
Average activation obtained = 2.8416999999999986
```

```
Average activation obtained = 0.47869999999999996
```

Under the obtained optimal control policy : AvgA = 2.84

Under no control policy control policy : AvgA =0.48

When p=0.05, the optimal control policy $\pi^*=[a^3,a^3,a^3,a^3,a^3,a^3,a^3,a^3,a^4,a^3,a^3,a^3,a^5,a^3,a^3,a^3]^T$

Compare to the AvgA under no control policy, the AvgA optimal control policy is greater, which is

## Part B

Repeat part(a) for p = 0.2 and p = 0.45.What is your observation? Compare the obtained optimal policy and AvgA for two noise parameters with the results of part(a).

p = 0.2

```
values= [[215.47953129 215.47953129 218.20647908 218.20647908 221.65428808
   221.65428808 221.65428808 221.65428808 218.72255113 215.47953129
   215.47953129 215.47953129 221.65428808 218.72255113 221.65428808
   221.65428808]]
policy= [2, 2, 2, 2, 2, 2, 2, 2, 3, 2, 2, 2, 4, 2, 2, 2]
```

```
Average activation obtained = 2.3330499999999996
```

```
Average activation obtained = 1.26795000000000004
```

When p = 0.2, the optimal control policy $\pi^*=[a^3,a^3,a^3,a^3,a^3,a^3,a^3,a^3,a^4,a^3,a^3,a^3,a^5,a^3,a^3,a^3]^T$
Under the obtained optimal control policy : AvgA =2.33
Under no control policy control policy : AvgA =1.26

p = 0.45

```
values= [[190.05579393 190.05579393 190.54235058 190.54235058 191.05474104
   191.05474104 191.05474104 191.05474104 190.62775028 190.05579393
   190.05579393 190.05579393 191.13724578 190.56574015 191.05474104
   191.05474104]]
[0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0]
```

```
Average activation obtained = 1.9083500000000004
```

```
Average activation obtained = 1.9124000000000008
```

When p = 0.45, the optimal control policy $\pi^*=[a^1,a^1,a^1,a^1,a^1,a^1,a^1,a^1,a^1,a^1,a^1,a^1,a^1,a^1,a^1,a^1]^T$
·Under the obtained optimal control policy : AvgA = 1.91
Under no control policy control policy : AvgA = 1.91

When p=0.2 or p= 0.42, the Average Activation is lower than part (a), since with the increase of probability, the noise is also increase,it might result in the selection of random action increase, so in this situation, the agent can hardly make choice to gain the greater reward.

## Part C

Similar to part(a) use p=0.05,an action $a^1$ in all states as the initial policy. Perform matrix- form

Policy Iteration method for computing the optimal control policy $\pi^*$. Compare the results with part (a). Are the obtained policies the same?

```
value= [[255.63715164 255.63715164 259.95441119 259.95441119 264.80249361
  264.80249361 264.80249361 264.80249361 260.31211153 255.63715164
  255.63715164 255.63715164 264.80249361 260.31211153 264.80249361
  264.80249361]]
policy= [2, 2, 2, 2, 2, 2, 2, 2, 3, 2, 2, 2, 4, 2, 2, 2]
```

```
Average activation obtained = 2.84065000000000006
```

The optimal control policy $\pi^*=[a^3,a^3,a^3,a^3,a^3,a^3,a^3,a^3,a^4,a^3,a^3,a^3,a^5,a^3,a^3,a^3]^T$
·Under the obtained optimal control policy : AvgA = 2.84

When p=0.05,using value iteration:

```
values= [[255.45551755 255.45551755 259.7727771  259.7727771  264.62085953
  264.62085953 264.62085953 264.62085953 260.13047744 255.45551755
  255.45551755 255.45551755 264.62085953 260.13047744 264.62085953
  264.62085953]]
policy= [2, 2, 2, 2, 2, 2, 2, 2, 3, 2, 2, 2, 4, 2, 2, 2]
```

```
Average activation obtained = 2.8416999999999986
```

The optimal control policy $\pi^*=[a^3,a^3,a^3,a^3,a^3,a^3,a^3,a^3,a^4,a^3,a^3,a^3,a^5,a^3,a^3,a^3]^T$
·Under the obtained optimal control policy : AvgA = 2.84

 The Policy iteration has the same policies with the Value Iteration when p=0.05. Compare the iterations with Policy iteration and Value iteration, the Policy Iteration has the less steps.