



Project #2

Problem 1

Consider a maze shown on the second page. This maze consists of several walls that the agent cannot enter and bumps and oils that moving to them have negative rewards. For simplicity, consider an 18×18 matrix, where each element is associated with one of the following:

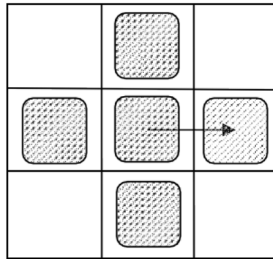
- *Empty*
- *Full (Wall)*
- *Bump*
- *Oil*

State Space (S): The state-space contains all cells in the maze except the walls, where the agent can possibly be there ($18 \times 18 - 76(\text{walls}) = 248$).

Action Space (A): The agent can take one of the four possible actions at any given state: *up (U)*, *down (D)*, *right (R)*, and *left (L)*.

Transition Probabilities: After choosing an action, the agent will either move to one of the neighborhood cells or stay in its current cell. After taking any action, with a probability of $1-p$, the agent moves to the anticipated state and, with an equal probability of $p/3$, will move to one of the other neighboring cells.

Consider the following example:

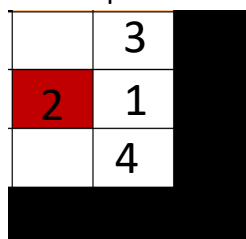


If action *right (R)* is selected, the agent: $\left\{ \begin{array}{l} \text{moves to the state in right with probability } 1 - p \\ \text{moves to the state in left with probability } p/3 \\ \text{moves to the the state in up with probability } p/3 \\ \text{moves to the state in down with probability } p/3 \end{array} \right.$

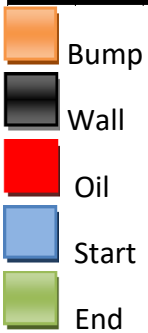
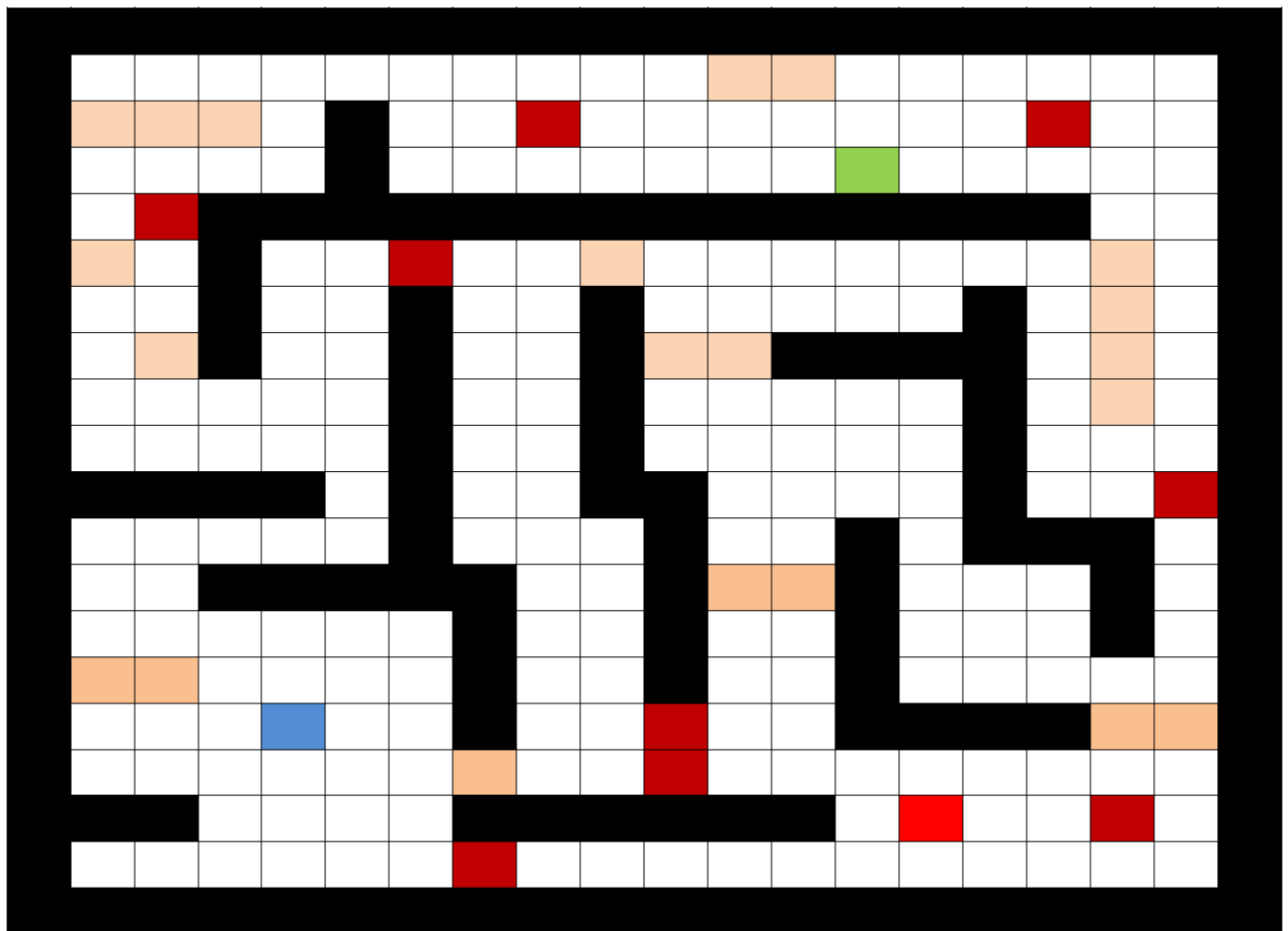
Notice that if any of the neighboring cells are wall, the agent stays in the current cell.

Reward Function:

The primary objective is to find the optimal policy which leads the agent to the goal state with the maximum accumulated reward. Therefore, we define the reward of taking any action -1 (delay), the reward for ending up to the oil states -5, the reward for ending up the bump states -10, and the reward for the goal state +200. For instance, this can be expressed for the following example as:



$$\left\{ \begin{array}{ll} R(1, L, 2) = -6 & - 5 \text{ for oil and } - 1 \text{ for taking an action} \\ R(1, L, 3) = -1 & \text{for taking an action} \\ R(1, L, 4) = -1 & \text{for taking an action} \\ R(1, L, 1) = -1 & \text{for taking an action} \end{array} \right.$$



1- 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. Solve the problem for the following cases:

- $p = 0.02, \gamma = 0.95, \theta = 0.01$ (Base Scenario)
- $p = 0.5, \gamma = 0.95, \theta = 0.01$ (Large Stochasticity Scenario)
- $p = 0.02, \gamma = 0.55, \theta = 0.01$ (Small Discount Factor Scenario)

Report the optimal value functions, the optimal policy, and the optimal path from start to end. Meanwhile, compare the results obtained in three parts (see the next page for an example of results that need to be provided).

2- Repeat the same process as part 1 with the vector-form Value Iteration method. Set initial state values to zero. Compare the results obtained from Value Iteration and Policy Iteration techniques.

Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$

1. Initialization

$V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation

Loop:

$\Delta \leftarrow 0$

Loop for each $s \in \mathcal{S}$:

$v \leftarrow V(s)$

$V(s) \leftarrow \sum_{s',r} p(s',r|s,\pi(s)) [r + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)

3. Policy Improvement

policy-stable \leftarrow *true*

For each $s \in \mathcal{S}$:

old-action $\leftarrow \pi(s)$

$\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a) [r + \gamma V(s')]$

If *old-action* $\neq \pi(s)$, then *policy-stable* \leftarrow *false*

If *policy-stable*, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

Value Iteration, for estimating $\pi \approx \pi_*$

Algorithm parameter: a small threshold $\theta > 0$ determining accuracy of estimation

Initialize $V(s)$, for all $s \in \mathcal{S}^+$, arbitrarily except that $V(\text{terminal}) = 0$

Loop:

$\Delta \leftarrow 0$

 Loop for each $s \in \mathcal{S}$:

$v \leftarrow V(s)$

$V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$

Output a deterministic policy, $\pi \approx \pi_*$, such that

$\pi(s) = \operatorname{argmax}_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$

Example of Results: For each combination of method and scenario, you need to provide the following set of plots:

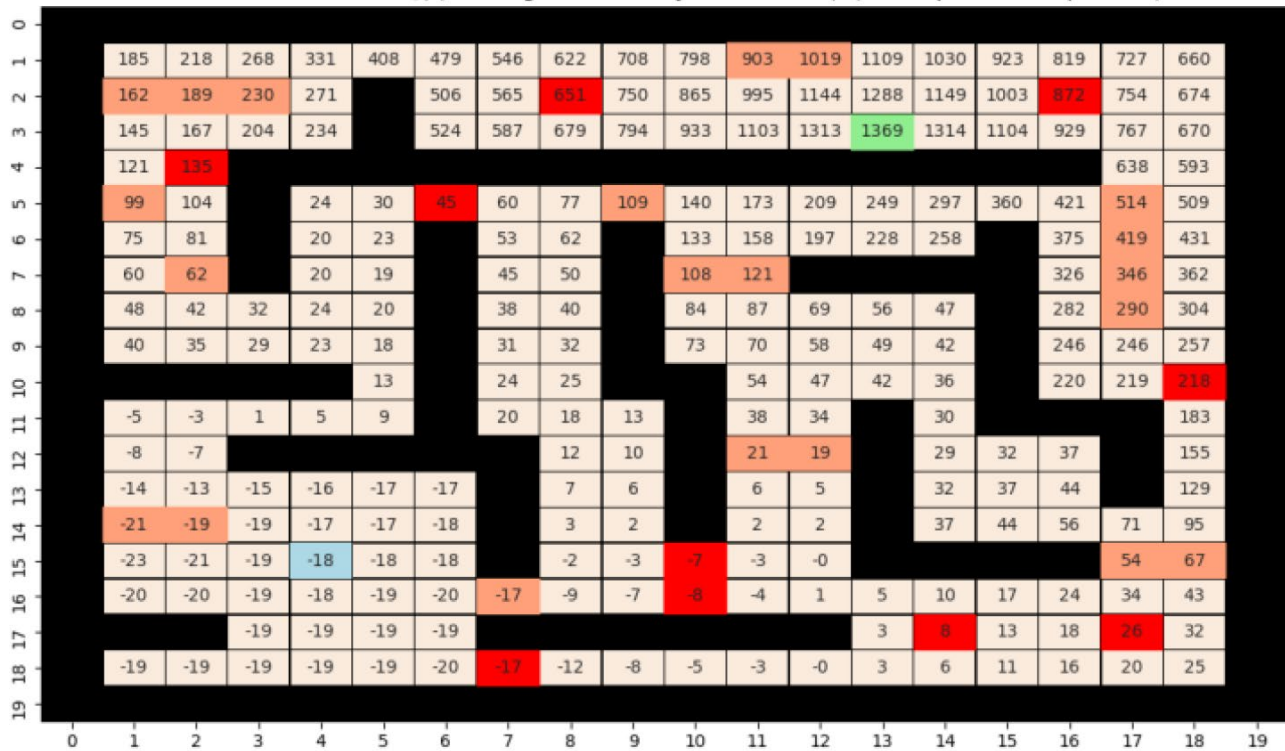


Figure 1 The optimal value function values at all states for the base scenario under the Policy Iteration method.

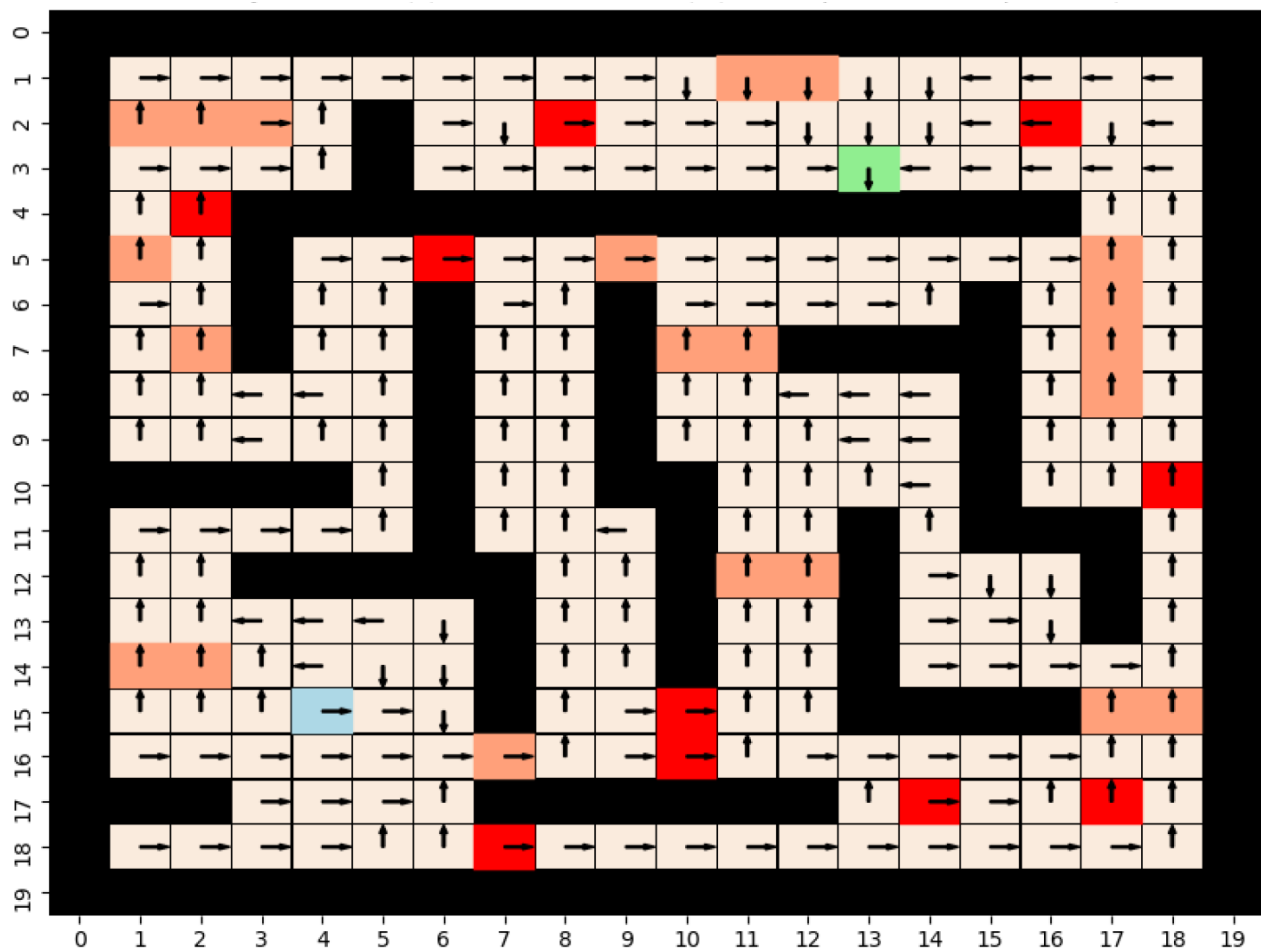


Figure 2 The optimal policy for all states for the base scenario under the Policy Iteration method.

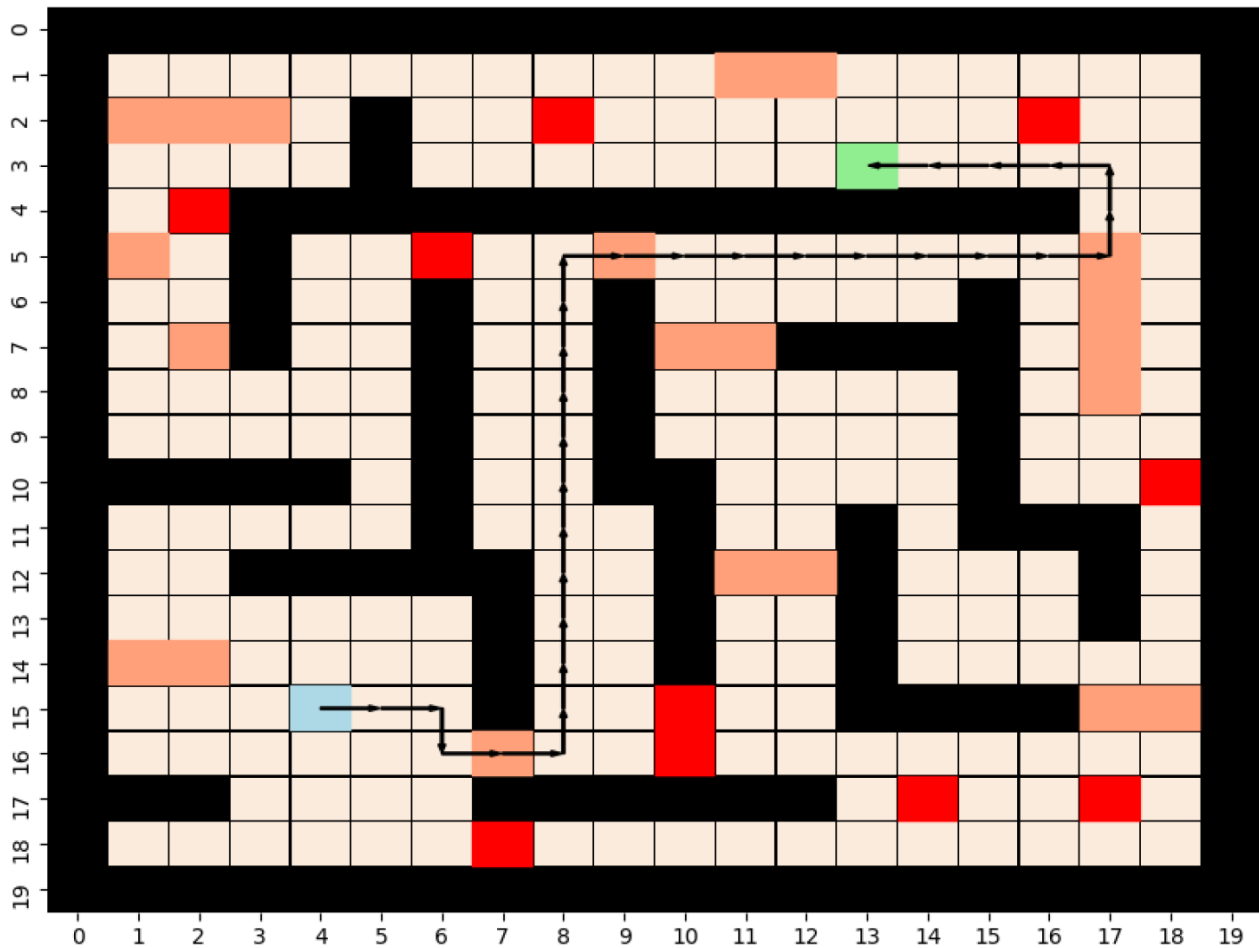


Figure 3 The optimal path from the start state to the end state for the base scenario under the Policy Iteration method.

You also need to compare the results obtained in all three scenarios and explain the difference and your thoughts on the obtained results. Please also report the number of iterations of Policy Iteration method (i.e., number of times that policy improvement methods has been performed) and the number of iteration of Value Iteration method.

Problem 2

The left plot of Fig. 1 represents the pathway diagram for a well-known biological network (the p53-Mdm2 negative feedback loop network). The network consists of 4 components/genes, represented in the state vector $\mathbf{s}_k = [\text{ATM}, \text{p53}, \text{Wpi}, \text{MDM2}]^T$.

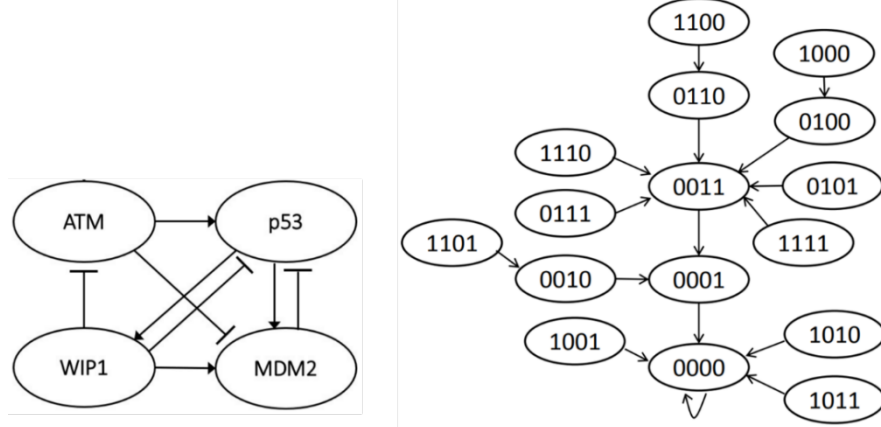


Figure 1 Pathway diagram for the p53-MDM2 network (left), and state transition diagram (right).

At any given time, the state value of each gene can be one of the following two values: 0 (OFF) and 1 (ON). Thus, the state space consists of $2^4 = 16$ states described below:

$\mathbf{s}^1 = [0,0,0,0]^T, \mathbf{s}^2 = [0,0,0,1]^T, \mathbf{s}^3 = [0,0,1,0]^T, \mathbf{s}^4 = [0,0,1,1]^T, \dots, \mathbf{s}^{15} = [1,1,1,0]^T, \mathbf{s}^{16} = [1,1,1,1]^T$. The action space is $A = \{\mathbf{a}^1 = [0,0,0,0]^T, \mathbf{a}^2 = [1,0,0,0]^T, \mathbf{a}^3 = [0,1,0,0]^T, \mathbf{a}^4 = [0,0,1,0]^T, \mathbf{a}^5 = [0,0,0,1]^T\}$, meaning that a single action at any time should be selected from the action space, i.e., $\mathbf{a}_{k-1} \in A = \{\mathbf{a}^1, \mathbf{a}^2, \mathbf{a}^3, \mathbf{a}^4, \mathbf{a}^5\}$.

The dynamics (state transition probabilities) of this system is governed by the following state transition:

$$\mathbf{s}_k = \underbrace{\begin{bmatrix} 0 & 0 & -1 & 0 \\ +1 & 0 & -1 & -1 \\ 0 & +1 & 0 & 0 \\ -1 & +1 & +1 & 0 \end{bmatrix}}_{\text{C: Connectivity Matrix}} \mathbf{s}_{k-1} \oplus \mathbf{a}_{k-1} \oplus \mathbf{n}_k \quad \text{Eq. (1)}$$

where \bar{v} maps the element of vector greater than 0 to 1, and smaller or equal to zero to 0, \mathbf{a}_{k-1} is a action input and \mathbf{n}_k is the state transition noise. The module-2 addition " \oplus " is like a component-wise XOR, such that

$$\begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix} \oplus \begin{bmatrix} 1 \\ 0 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \end{bmatrix}.$$

We assume that four elements of \mathbf{n}_k are identically distributed through a Bernoulli noise with parameter p . This can be expressed as: $\mathbf{n}_k(j) \sim \begin{cases} 1 & \text{with prob } p \\ 0 & \text{with prob } 1 - p \end{cases}$, for $j = 1, \dots, 4$.

The transition matrices for this system, which represent moving from any state to other states under different control inputs, are 16×16 matrices. The elements in the i th row and j th column of this matrix associated with action \mathbf{a} can be represented by

$$(M(\mathbf{a}))_{ij} = p^{\|\mathbf{s}^j - (\bar{C}\mathbf{s}^i \oplus \mathbf{a})\|_1} \times (1 - p)^{4 - \|\mathbf{s}^j - (\bar{C}\mathbf{s}^i \oplus \mathbf{a})\|_1}$$

where again $\|\mathbf{v}\|_1 = \sum_i |\mathbf{v}(i)|$ sums the absolute value of the elements of the vector \mathbf{v} .

For the system under no control, $\mathbf{a} = [0,0,0,0]^T$, the state transition diagram (without noise) is shown in the right plot of Fig. 1. It can be seen that the system without any process noise (deterministic case) ends up at state "0000" after a few transitions, meaning that the system spends most of its time in that state. The objective is to control/perturb the system so that the state values of all components become

1 (ON) at all time steps (i.e., continuing task). In other words, the system is desired to have as many active components as possible (ideally, spend most of the time on the state “1111”). This can be expressed through the following reward function:

$$R(s, \mathbf{a}, s') = 5 s'^{(1)} + 5 s'^{(2)} + 5 s'^{(3)} + 5 s'^{(4)} - |\mathbf{a}|$$

where $|\mathbf{a}|$ sums the absolute value of the elements of the vector \mathbf{a} . According to this reward function, the activation of each gene has the reward +5 and actions \mathbf{a}^2 to \mathbf{a}^5 have reward of -1.

Dynamic Programming: After constructing the controlled transition matrices, you need to use matrix-form Dynamic Programming with $\gamma = 0.95$.

- For $p = 0.05$, use matrix-form Value Iteration method for computing the optimal control policy π^* , 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 average activation of genes, $AvgA$ (definition is provided below), under the obtained optimal policy with $AvgA$ under no control policy (taking \mathbf{a}^1 in all states).
- 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).
- Similar to part (a), use $p = 0.05$, and action \mathbf{a}^1 in all states as the initial policy. Perform matrix-form Policy Iteration method for computing the optimal control policy π^* . Compare the results with part (a). Are the obtained policies the same?

- ❖ Please also report the number of iterations of Policy Iteration method (i.e., number of times that policy improvement methods has been performed) and the number of iteration of Value Iteration method.
- ❖ In addition to the final value function and policy (similar to problem 1), you need to evaluate the performance of the final obtained policy during the test in all parts. In episode i , start with a random initial state $s_0 \in \{s^1, \dots, s^{16}\}$; perform the action prescribed by the optimal policy, $\mathbf{a}_0 = \pi^*(s_0)$; use the transition matrix associated with \mathbf{a}_0 to compute s_1 (s_1 is a sample from $p(s'|s_0, \pi^*(s_0))$). Continue this to get the trajectory $D^i = \{(s_0, \mathbf{a}_0 = \pi^*(s_0), s_1, \mathbf{a}_1 = \pi^*(s_1), s_2, \dots, \mathbf{a}_{199} = \pi^*(s_{199}), s_{200})\}$. The average activation rate in episode i is $A^i = \frac{1}{200} \sum_{k=1}^{200} ||s_k||_1$, where the $0 \leq A^i \leq 4$, where the closer to 4 corresponds to perfect activation of all four genes at all times, whereas the values close to 1 denotes that a single gene (among 4) has been activated on average. Repeating the same process for 100 episodes, the average activation rate can be computed as:

$$AvgA = \frac{1}{100} \sum_{i=1}^{100} A^i$$

Questions about the project should be directed to TA, Mohammad Alali, at alali.m@northeastern.edu.

Value Iteration - Matrix Form

Transition matrices $M(a)$, $a \in A$, Reward $R_{ss'}^a$, Number of states N .

- Reward: $R_s^a = (M(a) \odot R_{ss'}^a) \mathbf{1}_{N \times 1}$, for $a \in A$.
 (Note: \odot is Hadamard Product)

- $V_0 = 0_{N \times 1}$, $k=0$

Repeat

- Value Iteration Backup: $V_{k+1} = \max_{a \in A} \left\{ R_s^a + \gamma M(a) V_k \right\}$
 (Note: \max is Row-wise)

- $k = k+1$

Until $\max_{i \in \{1, \dots, N\}} |V_k(i) - V_{k-1}(i)| < \theta$

- Optimal state values: $V^* = V_k$

- Optimal policy: $\pi^* = \arg \max_{a \in A} R_s^a + \gamma M(a) V^*$
 (Note: $\arg \max$ is Row-wise)

Policy Iteration: Matrix Form

Initialization: arbitrary Policy π , transition matrices $M(a)$, Reward $R_{ss'}^a$

- Reward: $R_s^a = (M(a) \odot R_{ss'}^a) \mathbf{1}_{N \times 1}$, for $a \in A$.
 (Note: \odot is Hadamard Product)

Policy Evaluation

- $M(\pi) = \begin{bmatrix} \text{Row 1 of } M(\pi(s^1)) \\ \vdots \\ \text{Row } N \text{ of } M(\pi(s^N)) \end{bmatrix}$, $R_s^\pi = \begin{bmatrix} \text{Element 1 of } R_s^{\pi(s^1)} \\ \vdots \\ \text{Element } N \text{ of } R_s^{\pi(s^N)} \end{bmatrix}$

- $V^\pi = (I - \gamma M(\pi))^{-1} R_s^\pi$

Policy Improvement

- $\pi' = \arg \max_{a \in A} R_s^a + \gamma M(a) V^\pi$
 (Note: $\arg \max$ is Row-wise)

If $\pi \neq \pi'$, $\pi = \pi'$ and go back to policy Evaluation step.

Else

- Optimal policy: $\pi^* = \pi'$

- Optimal state values: $V^* = V^{\pi^*}$