

# ECE64700: Homework IV

Bharath Keshavamurthy

## I. MARKOV DECISION PROCESSES - AN ENERGY HARVESTING RADIO NODE

### A. System Model

We consider a single energy harvesting radio node with an MDP agent running as a part of its decision engine solving for the optimal transmission policy considering channel fading states and battery energy levels.

The channel fading states  $g_k$  realize from the following set i.i.d over time with a uniform probability distribution  $\mathbb{P}(g_k = \gamma_i) = \theta_i = \frac{1}{N}$ .

$$\mathcal{G} \equiv \{\gamma_i, i = 1, 2, 3, \dots, N\} \quad (1)$$

The battery states  $Q_k$  realize from the following set.

$$\mathcal{Q} \equiv \{mE_0, m = 0, 1, 2, 3, \dots, M\} \quad (2)$$

The battery state update procedure is as follows.

$$Q_{k+1} = \min\{Q_k - w_k + a_k, ME_0\}$$

where,

$w_k \in \{0, E_0\}$  is the amount of energy spent by the radio node during transmission, and  $a_k \in \{0, E_0\}$  is modelled as a Bernoulli i.i.d process with  $\mathbb{P}(a_k = E_0) = p$  as the success probability and  $\mathbb{P}(a_k = 0) = (1 - p)$  as the failure probability.

**1) The State Space:** The state space for the MDP agent is denoted by  $\mathcal{S}$ .

In time-slot  $k$ , we represent the state of the process as a 2-tuple or a pair  $(Q_k, g_k)$  where,  $Q_k \in \mathcal{Q}$  is defined as the current battery level or the amount of energy available in the battery in time-slot  $k$ , and

$g_k \in \mathcal{G}$  is defined as the current state of the fading channel.

The dimensionality of the state space is given by  $|\mathcal{S}| = MN = 110$  possible states.

Therefore, concretely,

$$\mathcal{S} \equiv \{(Q_k, g_k) \mid Q_k \in \mathcal{Q}, g_k \in \mathcal{G}\} \quad (3)$$

**2) The Action Space:** The action space of the MDP agent is denoted by  $\mathcal{U}$ .

In time-slot  $k$ , the actions available to the MDP agent are,

$$u_k = \begin{cases} 0, & \text{if } Q_k = 0 \text{ and } g_k \in \mathcal{G} \\ \{1, 2, 3, \dots, i\}, & \text{if } Q_k \neq 0 \text{ and } g_k = \gamma_i \in \mathcal{G} \end{cases} \quad (4)$$

If the battery is dead, the radio node cannot transmit anything. Hence, the only available action in states with  $Q_k = 0$  is to remain idle, i.e.  $u_k = 0$ .

If there are units of energy in the battery (the units are discretized in steps of  $E_0$  as it is evident from equation (2)), the actions of the radio node depends on the realization of the fading state of the channel.

**3) The Transition Model:** The transition model of the process is given by,

$$(Q_{k+1}, g_{k+1}) = \begin{cases} (Q_k, g_{k+1}), & \text{if } u_k = 0 \text{ and } a_k = 0 \text{ with probability } \frac{1-p}{N} \\ (Q_k - E_0, g_{k+1}), & \text{if } u_k \neq 0 \text{ and } a_k = 0 \text{ with probability } \frac{1-p}{N} \\ (\min\{Q_k + E_0, ME_0\}, g_{k+1}), & \text{if } u_k = 0 \text{ and } a_k = E_0 \text{ with probability } \frac{p}{N} \\ (Q_k, g_{k+1}), & \text{if } u_k \neq 0 \text{ and } a_k = E_0 \text{ with probability } \frac{p}{N} \end{cases} \quad (5)$$

**4) The Reward Metric:** The reward metric is denoted by  $r_k = R(s_k, u_k)$ .

In time-slot  $k$ , we model the reward to the MDP agent as the number of bits it could transmit over the channel given the process state  $s_k = (Q_k, g_k) \in \mathcal{S}$ , i.e.  $r_k = u_k$ .

### B. The Problem Formulation

In order to solve for the optimal policy to be followed by the MDP agent residing in the decision engine of the energy harvesting radio node, we formulate the following problem which turns out to be the **Bellman Optimality Equation for Infinite Horizon Discounted Reward Problems**.

$$V^*(s) = \max_{u \in \mathcal{U}(s)} R(s, u) + \rho \sum_{s' \in \mathcal{S}} P_{s'|s,u} V^*(s') \quad (6)$$

where,

$\rho = 0.9$  is the discount factor,

$P_{s'|s,u}$  is the transition probability detailed in the transition model (5), and

$R(s, u)$  is the reward received in state  $s$  after taking action  $u$ .

### C. Solution Methodology and Performance

**1) The Algorithms:** We employ the Value Iteration Algorithm and the Policy Iteration Algorithm to solve for the optimal policy as shown below.

- **Value Iteration**

---

**Algorithm 1** The Value Iteration Algorithm for the given MDP agent

---

**Result:** The optimal policy  $\mu^*$  and optimal value function  $V^*(s), \forall s \in \mathcal{S}$

---

**Initialization:** Start with an initialization of the value function,  $V_0(s) = 10^{-5}, \forall s \in \mathcal{S}$ ;

**while**  $\forall s \in \mathcal{S}, V_{k+1}(s) \neq V_k(s)$  **do**

$$\left| \begin{array}{l} V_{k+1}(s) = \max_{u \in \mathcal{U}(s)} \left[ u + \rho \sum_{s' \in \mathcal{S}} P_{s'|s,u} V_k(s') \right] \\ u_k(s) = \arg \max_{u \in \mathcal{U}(s)} \left[ u + \rho \sum_{s' \in \mathcal{S}} P_{s'|s,u} V_k(s') \right] \end{array} \right|$$

**end**

Return  $V^*(s), \forall s \in \mathcal{S}$  and  $\mu^* = \{u^*(s) \mid s \in \mathcal{S}\}$  after convergence

---

- **Policy Iteration**

---

**Algorithm 2** The Policy Iteration Algorithm for the given MDP agent

---

**Result:** The optimal policy  $\mu^*$  and optimal value function  $V^*(s), \forall s \in \mathcal{S}$

---

**Initialization:** Start with an arbitrary policy  $\mu_0 = \{0 \text{ (idle)}, \forall s \in \mathcal{S}\}$  and an initialization of the value function,  $V_0(s) = 10^{-5}, \forall s \in \mathcal{S}$ ;

**while**  $\forall s \in \mathcal{S}, \mu_{k+1}(s) \neq \mu_k(s)$  **do**

**Policy Evaluation:**

$$V_{\mu_k}(s) = \left[ \mu_k(s) + \rho \sum_{s' \in \mathcal{S}} P_{s'|s,u} V_{\mu_k}(s') \right]$$

**Policy Improvement:**

$$\mu_{k+1}(s) = \arg \max_{u \in \mathcal{U}(s)} \left[ u + \rho \sum_{s' \in \mathcal{S}} P_{s'|s,u} V_{\mu_k}(s') \right]$$

**end**

Return  $V^*(s), \forall s \in \mathcal{S}$  and  $\mu^*$  after convergence

---

Please refer to the provided source code (*EnergyHarvestingNode.py*) for more details.

## 2) Results of the Algorithms:

### • Value Iteration

The optimal policy from the Value Iteration Algorithm is given below.

'0\_1': 0, '0\_2': 0, '0\_3': 0, '0\_4': 0, '0\_5': 0, '0\_6': 0, '0\_7': 0, '0\_8': 0, '0\_9': 0, '0\_10': 0,  
 '1\_1': 0, '1\_2': 0, '1\_3': 0, '1\_4': 0, '1\_5': 0, '1\_6': 0, '1\_7': 7, '1\_8': 8, '1\_9': 9, '1\_10': 10,  
 '2\_1': 0, '2\_2': 0, '2\_3': 0, '2\_4': 0, '2\_5': 5, '2\_6': 6, '2\_7': 7, '2\_8': 8, '2\_9': 9, '2\_10': 10,  
 '3\_1': 0, '3\_2': 0, '3\_3': 0, '3\_4': 0, '3\_5': 5, '3\_6': 6, '3\_7': 7, '3\_8': 8, '3\_9': 9, '3\_10': 10,  
 '4\_1': 0, '4\_2': 0, '4\_3': 0, '4\_4': 4, '4\_5': 5, '4\_6': 6, '4\_7': 7, '4\_8': 8, '4\_9': 9, '4\_10': 10,  
 '5\_1': 0, '5\_2': 0, '5\_3': 3, '5\_4': 4, '5\_5': 5, '5\_6': 6, '5\_7': 7, '5\_8': 8, '5\_9': 9, '5\_10': 10,  
 '6\_1': 0, '6\_2': 0, '6\_3': 3, '6\_4': 4, '6\_5': 5, '6\_6': 6, '6\_7': 7, '6\_8': 8, '6\_9': 9, '6\_10': 10,  
 '7\_1': 0, '7\_2': 0, '7\_3': 3, '7\_4': 4, '7\_5': 5, '7\_6': 6, '7\_7': 7, '7\_8': 8, '7\_9': 9, '7\_10': 10,  
 '8\_1': 0, '8\_2': 2, '8\_3': 3, '8\_4': 4, '8\_5': 5, '8\_6': 6, '8\_7': 7, '8\_8': 8, '8\_9': 9, '8\_10': 10,  
 '9\_1': 0, '9\_2': 2, '9\_3': 3, '9\_4': 4, '9\_5': 5, '9\_6': 6, '9\_7': 7, '9\_8': 8, '9\_9': 9, '9\_10': 10,  
 '10\_1': 1, '10\_2': 2, '10\_3': 3, '10\_4': 4, '10\_5': 5, '10\_6': 6, '10\_7': 7, '10\_8': 8, '10\_9': 9, '10\_10': 10

Here, 0\_1 : 0 refers to the optimal action of remaining idle, i.e.  $u_k = 0$ , if the process state is  $(Q_k = 0, g_k = \gamma_1)$ .

Similarly, 10\_7 : 7 refers to the optimal action of sending the maximum allowed bits of 7 if the process state is  $(Q_k = 10, g_k = \gamma_7)$ .

Some sample state value functions are:

$(Q_k = 0, g_k = \gamma_1)$ : 6.484672466785633  
 $(Q_k = 2, g_k = \gamma_4)$ : 17.801564174937106  
 $(Q_k = 5, g_k = \gamma_2)$ : 28.273127961666923  
 $(Q_k = 10, g_k = \gamma_1)$ : 37.55420277283401  
 $(Q_k = 10, g_k = \gamma_{10})$ : 46.55420700996909

## • Policy Iteration

The optimal policy from the Policy Iteration algorithm is given by,

'0\_1': 0, '0\_2': 0, '0\_3': 0, '0\_4': 0, '0\_5': 0, '0\_6': 0, '0\_7': 0, '0\_8': 0, '0\_9': 0, '0\_10': 0,  
 '1\_1': 0, '1\_2': 0, '1\_3': 0, '1\_4': 0, '1\_5': 0, '1\_6': 0, '1\_7': 7, '1\_8': 8, '1\_9': 9, '1\_10': 10,  
 '2\_1': 0, '2\_2': 0, '2\_3': 0, '2\_4': 0, '2\_5': 5, '2\_6': 6, '2\_7': 7, '2\_8': 8, '2\_9': 9, '2\_10': 10,  
 '3\_1': 0, '3\_2': 0, '3\_3': 0, '3\_4': 0, '3\_5': 5, '3\_6': 6, '3\_7': 7, '3\_8': 8, '3\_9': 9, '3\_10': 10,  
 '4\_1': 0, '4\_2': 0, '4\_3': 0, '4\_4': 4, '4\_5': 5, '4\_6': 6, '4\_7': 7, '4\_8': 8, '4\_9': 9, '4\_10': 10,  
 '5\_1': 0, '5\_2': 0, '5\_3': 3, '5\_4': 4, '5\_5': 5, '5\_6': 6, '5\_7': 7, '5\_8': 8, '5\_9': 9, '5\_10': 10,  
 '6\_1': 0, '6\_2': 0, '6\_3': 3, '6\_4': 4, '6\_5': 5, '6\_6': 6, '6\_7': 7, '6\_8': 8, '6\_9': 9, '6\_10': 10,  
 '7\_1': 0, '7\_2': 0, '7\_3': 3, '7\_4': 4, '7\_5': 5, '7\_6': 6, '7\_7': 7, '7\_8': 8, '7\_9': 9, '7\_10': 10,  
 '8\_1': 0, '8\_2': 2, '8\_3': 3, '8\_4': 4, '8\_5': 5, '8\_6': 6, '8\_7': 7, '8\_8': 8, '8\_9': 9, '8\_10': 10,  
 '9\_1': 0, '9\_2': 2, '9\_3': 3, '9\_4': 4, '9\_5': 5, '9\_6': 6, '9\_7': 7, '9\_8': 8, '9\_9': 9, '9\_10': 10,  
 '10\_1': 1, '10\_2': 2, '10\_3': 3, '10\_4': 4, '10\_5': 5, '10\_6': 6, '10\_7': 7, '10\_8': 8, '10\_9': 9, '10\_10': 10

Here, 0\_1 : 0 refers to the optimal action of remaining idle, i.e.  $u_k = 0$ , if the process state is ( $Q_k = 0$ ,  $g_k = \gamma_1$ ).

Similarly, 10\_7 : 7 refers to the optimal action of sending the maximum allowed bits of 7 if the process state is ( $Q_k = 10$ ,  $g_k = \gamma_7$ ).

Some sample state value functions are:

( $Q_k = 0$ ,  $g_k = \gamma_1$ ): 6.484553119956839  
 ( $Q_k = 2$ ,  $g_k = \gamma_4$ ): 17.801464754820753  
 ( $Q_k = 5$ ,  $g_k = \gamma_2$ ): 28.27304768957006  
 ( $Q_k = 10$ ,  $g_k = \gamma_1$ ): 37.554155315117185  
 ( $Q_k = 10$ ,  $g_k = \gamma_{10}$ ): 46.554160324453505

Please refer to the provided Console Logs for more details on the results of these two algorithms.

*Although Value Iteration and Policy Iteration are comparable in their performance due to smaller state and action spaces in this process, we can see from the plots that Value Iteration is slightly faster than Policy Iteration.*

Convergence Visualization of the Value Iteration and Policy Iteration Algorithms for the Decision Engine in the given Energy Harvesting System considering system state ( $Q_k = 9, g_k = \gamma_{10}$ )

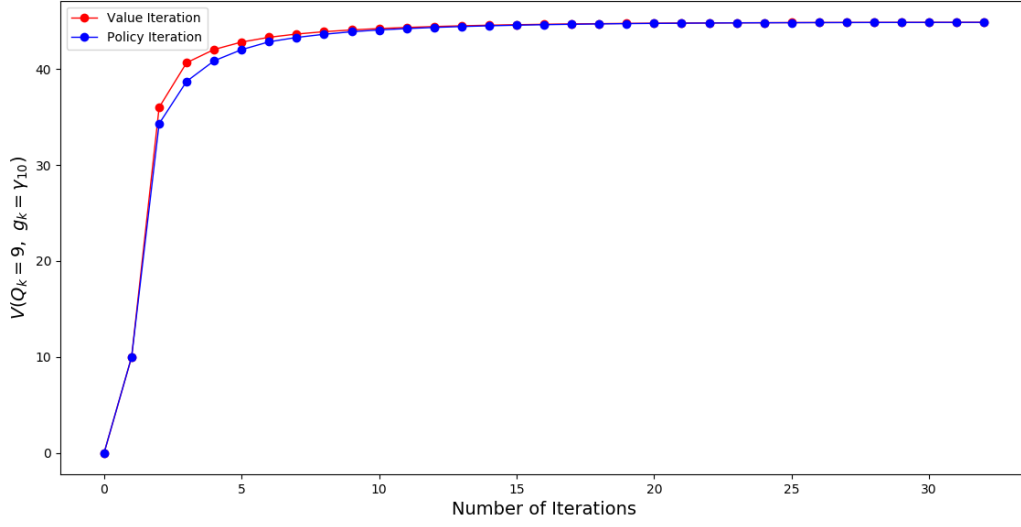


Fig. 1. Convergence Visualization of the Value Iteration and Policy Iteration Algorithms for the Decision Engine in the given Energy Harvesting System considering system state ( $Q_k = 9, g_k = \gamma_{10}$ )

Convergence Visualization of the Value Iteration and Policy Iteration Algorithms for the Decision Engine in the given Energy Harvesting System considering system state ( $Q_k = 9, g_k = \gamma_{10}$ )

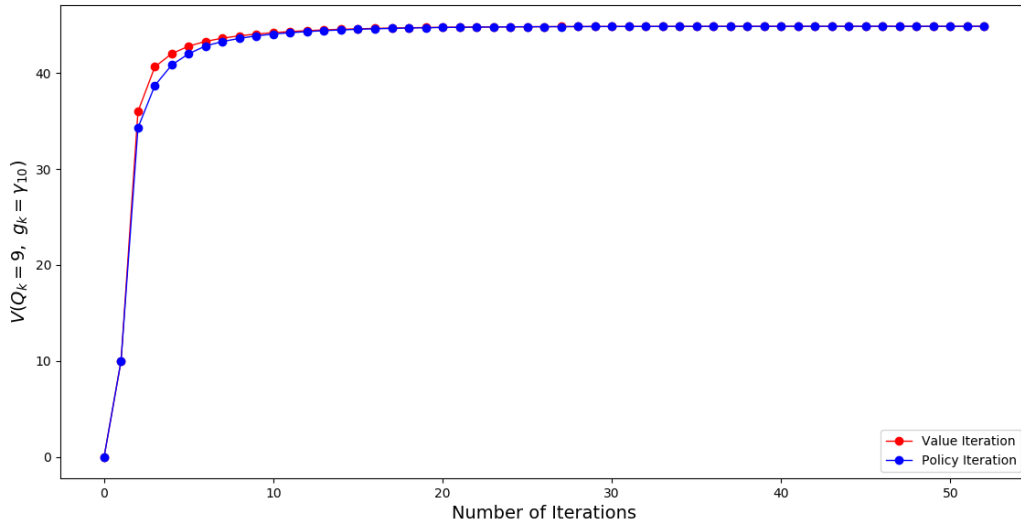


Fig. 2. Convergence Visualization of the Value Iteration and Policy Iteration Algorithms for the Decision Engine in the given Energy Harvesting System considering system state ( $Q_k = 9, g_k = \gamma_{10}$ ) - Increased number of iterations to prove convergence

Convergence Visualization of the Value Iteration and Policy Iteration Algorithms for the Decision Engine in the given Energy Harvesting System considering system state ( $Q_k = 7, g_k = \gamma_4$ )

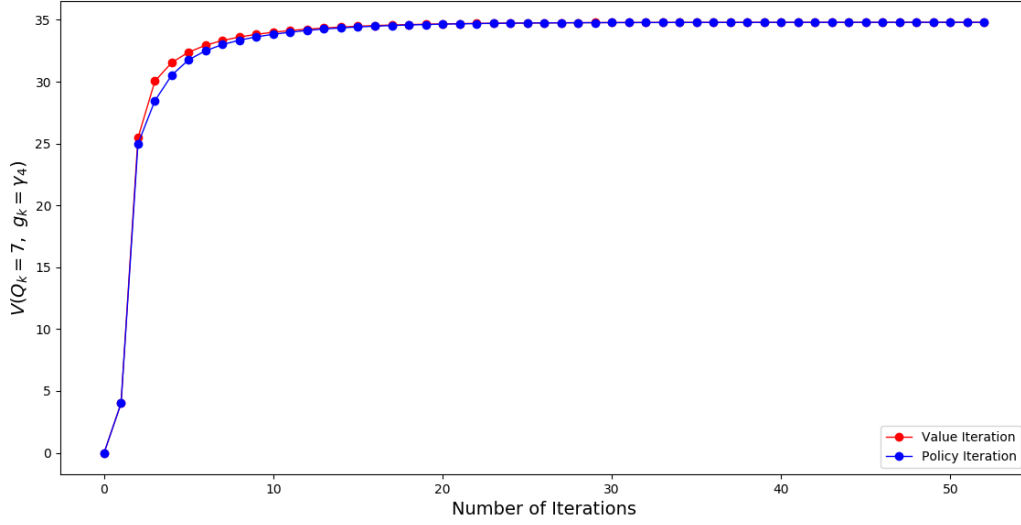


Fig. 3. Convergence Visualization of the Value Iteration and Policy Iteration Algorithms for the Decision Engine in the given Energy Harvesting System considering system state ( $Q_k = 7, g_k = \gamma_4$ )

Convergence Visualization of the Value Iteration and Policy Iteration Algorithms for the Decision Engine in the given Energy Harvesting System considering system state ( $Q_k = 3, g_k = \gamma_7$ )

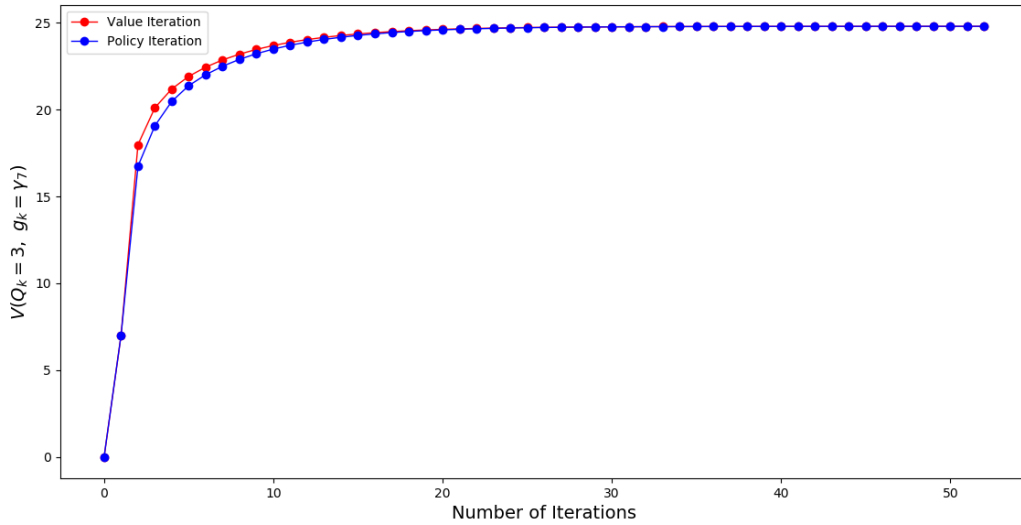


Fig. 4. Convergence Visualization of the Value Iteration and Policy Iteration Algorithms for the Decision Engine in the given Energy Harvesting System considering system state ( $Q_k = 3, g_k = \gamma_7$ )

Convergence Visualization of the Value Iteration and Policy Iteration Algorithms for the Decision Engine in the given Energy Harvesting System considering system state ( $Q_k = 1, g_k = \gamma_9$ )

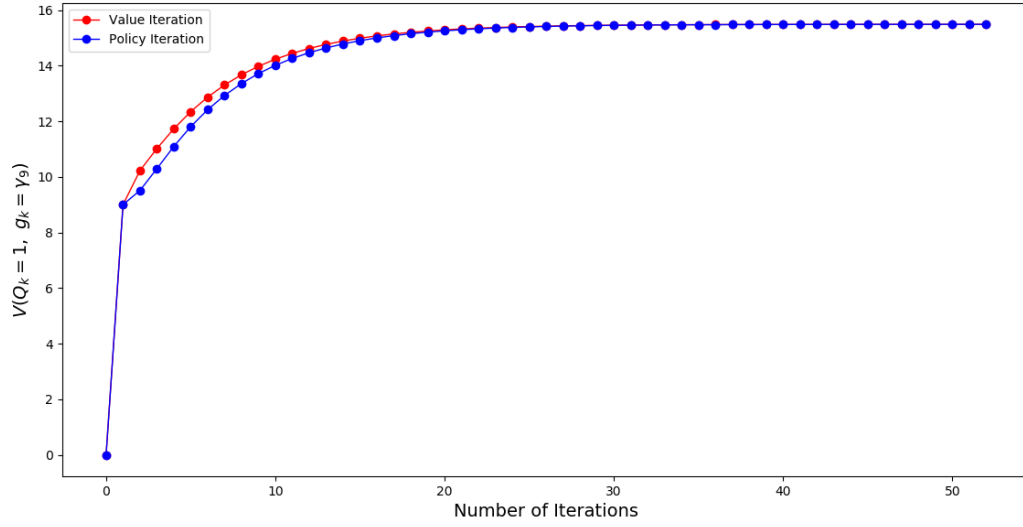


Fig. 5. Convergence Visualization of the Value Iteration and Policy Iteration Algorithms for the Decision Engine in the given Energy Harvesting System considering system state ( $Q_k = 1, g_k = \gamma_9$ )

Optimal policy versus Battery State and Fading State for the Decision Engine in the given Energy Harvesting System

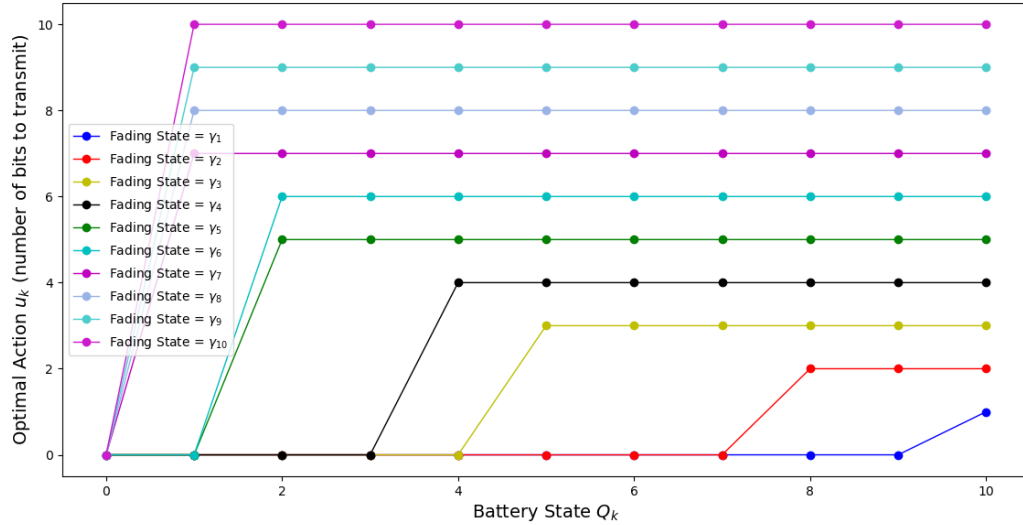


Fig. 6. Optimal policy versus Battery State and Fading State for the Decision Engine in the given Energy Harvesting System