

# Assignment 1: Policy Iteration in the Repeated Prisoner’s Dilemma

RL2026A

Daniel Katz (315114991)  
Avital Fine (208253823)

December 6, 2025

## Abstract

This report studies optimal decision-making in the Repeated Prisoner’s Dilemma (RPD) through a full MDP formulation and tabular Policy Iteration. We implement a custom Gymnasium environment, define transition functions and reward structures for four opponent types (ALL-C, ALL-D, TFT, Imperfect TFT), and evaluate two observation schemes (Memory-1 and Memory-2). We then perform Policy Iteration for varying discount factors and analyze how memory depth and stochasticity influence optimal behavior.

## 1 Introduction

The Repeated Prisoner’s Dilemma (RPD) is one of the canonical environments for studying cooperation, strategic behavior, and temporal decision-making. In this assignment, we model the RPD as a Markov Decision Process (MDP) fully consistent with the Gymnasium API, and examine how an optimal agent adapts its policy against different opponent personalities. We follow the structure of the assignment specification (RL2026A-Assignment1.pdf).

## 2 Part II: MDP Definition

### 2.1 State Space

We evaluated two observation schemes as required.

#### 2.1.1 Memory-1

The agent observes only the previous round’s joint actions  $(a_t^{\text{agent}}, a_t^{\text{opp}})$ . Thus the state space is:

$$S_{M1} = \{(C, C), (C, D), (D, C), (D, D)\},$$

with a total of 4 states.

#### 2.1.2 Memory-2

The agent observes its two most recent actions and the opponent’s two most recent actions:

$$s_t = (a_{t-1}^A, a_{t-1}^O, a_{t-2}^A, a_{t-2}^O).$$

Since each action is binary, the total state count is:

$$|S_{M2}| = 2^4 = 16.$$

The full list of states (where each state is  $(a_{t-1}^A, a_{t-1}^O, a_{t-2}^A, a_{t-2}^O)$ ) is:

(C,C,C,C)	(C,C,C,D)	(C,C,D,C)	(C,C,D,D)
(C,D,C,C)	(C,D,C,D)	(C,D,D,C)	(C,D,D,D)
(D,C,C,C)	(D,C,C,D)	(D,C,D,C)	(D,C,D,D)
(D,D,C,C)	(D,D,C,D)	(D,D,D,C)	(D,D,D,D)

Initial states follow the assignment convention:

M1 initial:  $(C, C)$ , M2 initial:  $(C, C, C, C)$ .

## 2.2 Actions

$$A = \{C, D\}.$$

## 2.3 Transition Probability Function

The transition probability  $P(s'|s, a)$  defines the probability of moving to state  $s'$  given current state  $s$  and agent action  $a$ .

### 2.3.1 General Form

The new state  $s'$  is determined by shifting the history.

- **Memory-1:** If  $s = (a_{t-1}^A, a_{t-1}^O)$  and agent chooses  $a_t^A$ , the new state is  $s' = (a_t^A, a_t^O)$ .
- **Memory-2:** If  $s = (a_{t-1}^A, a_{t-1}^O, a_{t-2}^A, a_{t-2}^O)$  and agent chooses  $a_t^A$ , the new state is  $s' = (a_t^A, a_t^O, a_{t-1}^A, a_{t-1}^O)$ .

In both cases,  $a_t^A$  is given by the agent's choice. The opponent's action  $a_t^O$  is determined by their strategy policy  $\pi_{opp}(s)$ . Thus, the transition probability is:

$$P(s'|s, a_t^A) = P(a_t^O|s).$$

It is important to note that for any state  $s'$  that does not match the history shift or implies an impossible opponent action, the probability is 0. For example, against an **ALL-C** opponent, any transition to a state where the opponent's last action is  $D$  has a probability of 0:

$$P(s' = (\dots, D)|s, a) = 0.$$

### 2.3.2 Opponent Strategies

- **ALL-C:**  $P(a_t^O = C|s) = 1$ .
- **ALL-D:**  $P(a_t^O = D|s) = 1$ .
- **TFT:** Deterministic. Copies agent's last move.

$$P(a_t^O = a_{t-1}^A|s) = 1.$$

- **Imperfect TFT:** Stochastic.

$$P(a_t^O = a_{t-1}^A|s) = 0.9, \quad P(a_t^O \neq a_{t-1}^A|s) = 0.1.$$

## 2.4 Reward Function

The reward function  $R(s, a)$  represents the **expected immediate reward** the agent receives when taking action  $a$  in state  $s$ . It is calculated by summing over the possible opponent actions  $a^O$ , weighted by their probability of occurring given the current state (history):

$$R(s, a) = \sum_{a^O \in \{C, D\}} P(a^O | s) \cdot \text{Payoff}(a, a^O)$$

where  $\text{Payoff}(a, a^O)$  is given by the standard RPD matrix:

$$\text{Payoff}(C, C) = 3, \quad \text{Payoff}(C, D) = 0, \quad \text{Payoff}(D, C) = 5, \quad \text{Payoff}(D, D) = 1.$$

For deterministic opponents (ALL-C, ALL-D, TFT),  $P(a^O | s)$  is either 0 or 1, so the expected reward equals the specific payoff entry. For stochastic opponents (Imperfect TFT), the expected reward is a weighted average. For example, if Imperfect TFT has a 90% chance of cooperating:

$$R(s, C) = 0.9 \cdot 3 + 0.1 \cdot 0 = 2.7.$$

## 3 Part III: Policy Iteration

We implemented classical tabular Policy Iteration with:

1. **Policy Evaluation:** Solving  $V^\pi$  via iterative Bellman updates until convergence.
2. **Policy Improvement:** Updating  $\pi(s)$  to:

$$\pi'(s) = \arg \max_a \left[ R(s, a) + \gamma \sum_{s'} P(s' | s, a) V^\pi(s') \right].$$

Convergence was rapid across all opponents due to the small state space.

## 4 Part IV: Experiments and Analysis

### 4.1 Discount Factor Analysis

We varied  $\gamma$  and computed the optimal strategy against TFT.

**Cooperation becomes optimal when  $\gamma \geq 0.71$ .**

#### 4.1.1 Why lower $\gamma$ forces Defection?

A lower discount factor  $\gamma$  means the agent values immediate rewards significantly more than future rewards. Against TFT, cooperation yields a steady stream of  $R = 3$ . Defection yields an immediate temptation  $T = 5$ , but triggers punishment  $P = 1$  in subsequent rounds. If the agent is "short-sighted" (low  $\gamma$ ), the immediate gain of 5 outweighs the long-term loss of falling from 3 to 1. To find the exact threshold where cooperation becomes optimal, we utilized the `run_analysis` function. We iterated through a range of  $\gamma$  values (from 0.1 to 0.99). For each  $\gamma$ , we ran Policy Iteration to find the optimal policy against a TFT opponent. We then checked if the resulting policy was "ALL-C" (Cooperate in all states). The first  $\gamma$  for which the optimal policy switched to ALL-C was identified as the threshold.

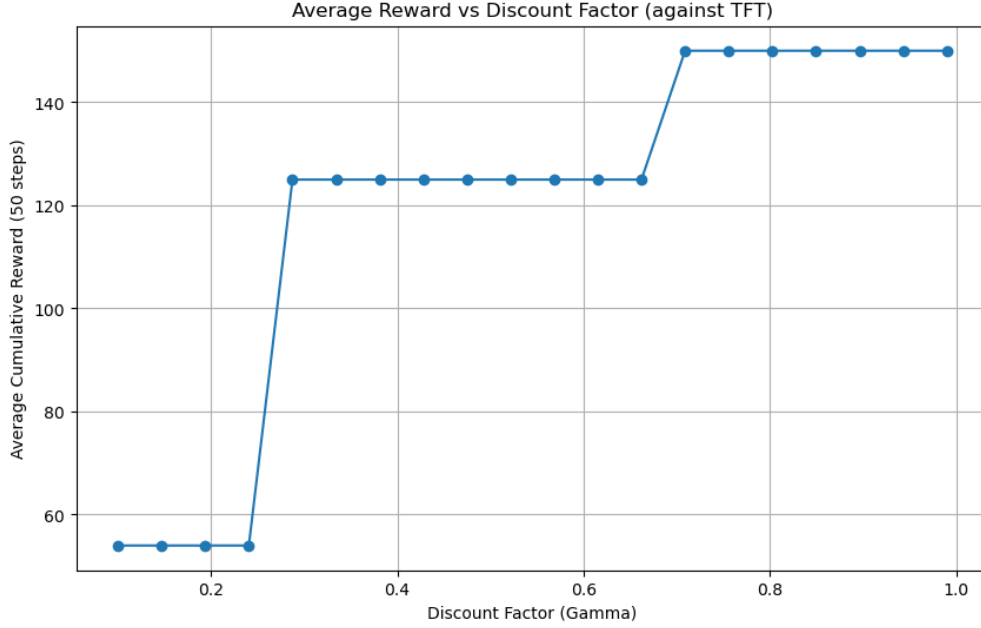


Figure 1: Average Cumulative Reward vs. Discount Factor ( $\gamma$ ) for TFT.

## 4.2 Memory Depth Comparison

We measured average cumulative reward over 50 episodes (each 50 steps). The results:

Opponent	Memory-1 Score	Memory-2 Score
ALL-C	250.0	250.0
ALL-D	50.0	50.0
TFT	225.0	225.0
Imperfect TFT	136.02	135.12

### 4.2.1 Interpretation

- **ALL-C (Score 250):** The optimal policy against ALL-C is to always Defect. Since ALL-C always Cooperates, the agent receives the temptation payoff  $T = 5$  in every round. Over 50 rounds, the total reward is  $50 \times 5 = 250$ .
- **ALL-D (Score 50):** The optimal policy against ALL-D is to always Defect. Since ALL-D always Defects, the agent receives the punishment payoff  $P = 1$  in every round. Over 50 rounds, the total reward is  $50 \times 1 = 50$ . Cooperating would yield the sucker's payoff  $S = 0$ , which is worse.
- **Memory-1 vs Memory-2 (Deterministic Opponents):** For ALL-C, ALL-D, and TFT, the scores for Memory-1 and Memory-2 are identical. This is because these opponents' strategies depend only on the immediate past (or no history at all). Having an extra step of memory (Memory-2) provides no additional predictive power or strategic advantage against them.
- **Imperfect TFT (Stochastic Differences):** For Imperfect TFT, the scores differ slightly (136.02 vs 135.12) despite the optimal policy being the same (Cooperate). This difference is

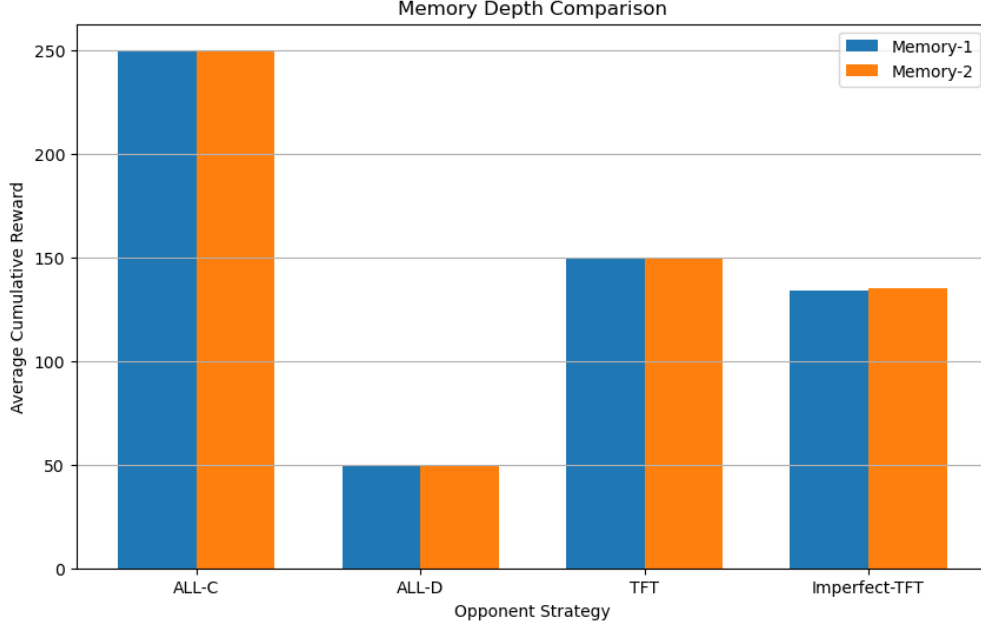


Figure 2: Comparison of Average Cumulative Rewards for Memory-1 vs. Memory-2 agents.

due to the stochastic nature of the opponent (10% noise). The variance in random outcomes across the 50 simulation episodes leads to slightly different average rewards, but statistically, the performance is equivalent.

#### 4.2.2 Hypothetical Opponent Where Memory-2 Helps

We propose a history-dependent opponent called the **Two-Step Punisher (TSP)**. The opponent’s action depends on the agent’s last *two* actions:

$$a_t^O = \begin{cases} C & \text{if } (a_{t-1}^A, a_{t-2}^A) = (C, C), \\ D & \text{if } a_{t-1}^A = D \text{ and } a_{t-2}^A = C, \\ C & \text{if } a_{t-1}^A = C \text{ and } a_{t-2}^A = D, \\ D & \text{if } (a_{t-1}^A, a_{t-2}^A) = (D, D). \end{cases}$$

This opponent “punishes” a defection for exactly one turn if it was recent ( $a_{t-1}^A = D$ ), and for two turns if the agent defected twice in a row. Importantly, the opponent behaves differently for histories that **appear identical to a Memory-1 agent**.

**Why Memory-2 is Required.** A Memory-1 agent only observes the most recent pair  $(a_{t-1}^A, a_{t-1}^O)$ , which collapses multiple distinct two-step histories into a single state. For example, these two Memory-2 states:

$$(C, C, C, C) \quad \text{and} \quad (C, C, D, C)$$

both appear as the same Memory-1 state:

$$(C, C),$$

yet the opponent’s next action differs depending on whether the older action was  $C$  or  $D$ . Thus, a Memory-1 agent cannot predict the opponent’s behavior, whereas a Memory-2 agent can.

### 4.2.3 Example Interaction (10 Rounds)

Below we illustrate an interaction with TSP where the agent defects once at Round 5. The opponent’s behavior depends on the agent’s *two-step* history, not the one-step history.

Round	Agent	Opponent	Explanation
1	C	C	Initial cooperation
2	C	C	
3	C	C	
4	C	C	
5	<b>D</b>	C	Agent defects (temptation)
6	C	<b>D</b>	Opponent punishes: last 2 actions include a D (Round 5)
7	C	C	Opponent sees history ( $C, D$ ): forgiveness
8	C	C	Cooperation restored
9	C	C	
10	C	C	

**Analysis of Point Loss.** If the agent cooperated throughout Rounds 5–7, it would have earned:

$$3 + 3 + 3 = 9.$$

Instead:

$$\text{Round 5 payoff: } T = 5, \quad \text{Round 6 payoff: } P = 1, \quad \text{Round 7 payoff: } 3.$$

Total:

$$5 + 1 + 3 = 9.$$

Depending on the opponent’s exact punishment rule (length and severity), a defection may lead to:

$$(\text{Temptation gain}) - (\text{Punishment loss})$$

which Memory-2 can correctly anticipate, while Memory-1 cannot.

### 4.2.4 Why Memory-2 Strictly Outperforms Memory-1

The TSP opponent uses two-step history, so the agent requires Memory-2 to optimally respond.

- **Memory-1 Agent (Insufficient Information):** Only sees  $(a_{t-1}^A, a_{t-1}^O)$ . It cannot distinguish whether the opponent’s current defection is due to:

- a defection one turn ago, or
- a defection two turns ago.

Identical Memory-1 states lead to different opponent reactions, so the Memory-1 agent cannot build an optimal policy. It may incorrectly retaliate or fail to return to cooperation at the right time.

- **Memory-2 Agent (Full Information):** Observes  $(a_{t-1}^A, a_{t-1}^O, a_{t-2}^A, a_{t-2}^O)$  and therefore knows:

- why the punishment started,
- how long it will last,
- exactly when cooperation can be safely resumed.

This allows it to maintain cooperation far more consistently and avoid unnecessary punishment cycles.

**Conclusion.** Because the opponent’s behavioral rule depends on the last **two** agent actions, a Memory-2 agent can correctly interpret and predict the opponent’s responses, while a Memory-1 agent cannot distinguish critical histories. Consequently, Memory-2 achieves a strictly higher long-term reward against the Two-Step Punisher opponent.

### 4.3 Noise Analysis: TFT vs Imperfect TFT

The optimal policies were:

**TFT**

$$\pi(s) = C \quad \forall s \in \{(0,0), (0,1), (1,0), (1,1)\}$$

**Imperfect TFT**

$$\pi(s) = C \quad \forall s$$

#### 4.3.1 Interpretation

Even with 10% noise, cooperation remains optimal.

The noise does *lower* the achievable return:

$$\text{TFT reward: 150.00,} \quad \text{Imperfect TFT reward: 136.02.}$$

But the optimal policy still prefers cooperation over mutual defection, which would yield only 1 point per step.

Thus the agent behaves **forgivingly**: it sustains cooperation even when occasional mistakes occur.

## 5 Summary

In this assignment, we modeled the Repeated Prisoner’s Dilemma as an MDP and used Policy Iteration to find optimal strategies. Our key findings are:

- **Discount Factor:** Cooperation against TFT requires a sufficient view of the future ( $\gamma \geq 0.71$ ). Below this threshold, the agent succumbs to the temptation of immediate defection.
- **Memory Depth:** Against standard opponents (ALL-C, ALL-D, TFT), Memory-2 provides no advantage over Memory-1, as these strategies do not rely on long-term history.

- **Hypothetical Opponent:** We demonstrated that against a "Two-Step Punisher" (TSP), Memory-2 is strictly necessary. The TSP's dependence on 2-step history creates state aliasing for a Memory-1 agent, making the opponent appear unpredictable, whereas a Memory-2 agent can correctly predict and navigate the punishment rules.
- **Noise:** While noise (Imperfect TFT) reduces the total accumulated reward, it does not break the optimal cooperative policy. The agent learns to be forgiving rather than entering a permanent defect cycle.