

# CSCI B659: Reinforcement Learning

## Assignment 1: Experimental Results

LJ Huang

Spring 2025

### Contents

<b>1 Task 1: Comparing VI, PI, and MPI</b>	<b>2</b>
1.1 Experimental Plots for Task 1 . . . . .	2
1.1.1 Value Iteration (VI) . . . . .	2
1.1.2 Policy Iteration (PI) . . . . .	2
1.1.3 Modified Policy Iteration (MPI) . . . . .	3
1.2 Observations on Task 1 . . . . .	3
<b>2 Task 2: Model-Based RL with Corrupted Transition Models</b>	<b>4</b>
2.1 Experimental Plot for Task 2 . . . . .	4
2.2 Observations on Task 2 . . . . .	5
<b>3 Task 3: Naive Model-Based RL (MBRL)</b>	<b>6</b>
3.1 Experimental Plot for Task 3 . . . . .	6
3.2 Observations on Task 3 . . . . .	6
<b>A Supplementary Terminal Output Results</b>	<b>7</b>
A.1 Value Iteration (VI) . . . . .	7
A.2 Policy Iteration (PI) . . . . .	8
A.3 Modified Policy Iteration (MPI) . . . . .	8
<b>B README File</b>	<b>9</b>

# 1 Task 1: Comparing VI, PI, and MPI

In this task, we implement Value Iteration (VI), Policy Iteration (PI), and Modified Policy Iteration (MPI) on the Frozen Lake environment. We use a discount factor  $\gamma = 0.999$  and a stopping gap of 0.001. The algorithms are evaluated by running the current policy for 500 episodes to compute the mean discounted return. For VI the policy is evaluated every 10 iterations; for PI and MPI the policy is evaluated after each outer iteration. We also track the computational effort by counting the number of single action backups.

## 1.1 Experimental Plots for Task 1

### 1.1.1 Value Iteration (VI)

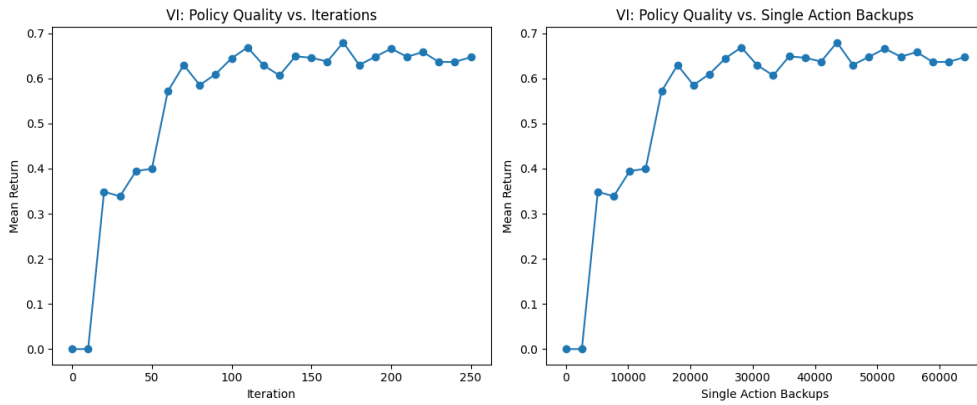


Figure 1: VI: Policy Quality vs. Iterations

### 1.1.2 Policy Iteration (PI)

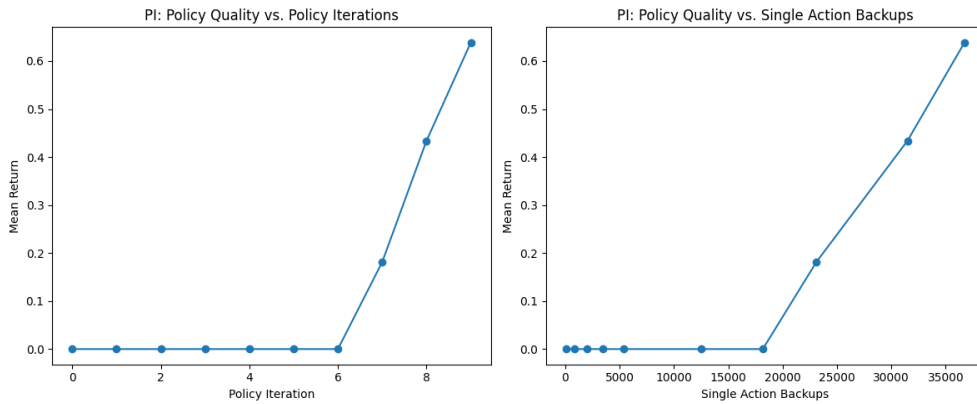


Figure 2: PI: Policy Quality vs. Policy Iterations

### 1.1.3 Modified Policy Iteration (MPI)

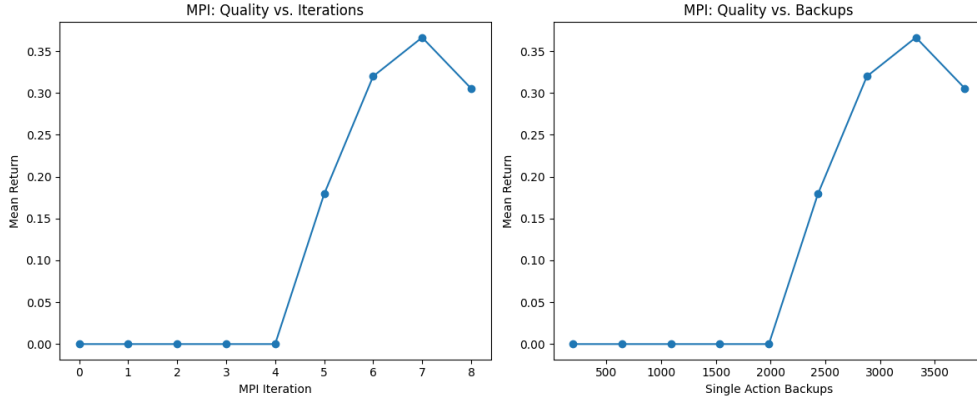


Figure 3: MPI: Policy Quality vs. MPI Iterations

## 1.2 Observations on Task 1

In this section, we analyze how the algorithms compare when considering the total computational effort. Here, computational effort is measured both in terms of the number of iterations and in terms of the total number of single action backups performed.

Our experimental results show that:

- **Policy Iteration (PI)** achieves the highest policy quality for a given amount of computational work.
- **Modified Policy Iteration (MPI)** requires more backups than PI but performs better than VI.
- **Value Iteration (VI)** requires the most computational effort (i.e., the largest number of backups) to reach a similar level of performance.

This ordering, **PI** > **MPI** > **VI**, indicates that, in our experiments, PI is the most computationally efficient method—providing higher mean returns per unit of computational effort—while VI is the least efficient.

## 2 Task 2: Model-Based RL with Corrupted Transition Models

In this task, we evaluate the robustness of planning by corrupting the true transition model. For each noise level  $\alpha$ , the model is distorted as follows:

1. For each state  $s$  and action  $a$ , let  $S'$  be the set of next states with  $p(s'|s, a) > 0$ , and let  $p = [p_1, p_2, \dots, p_k]$  be the original probability vector.
2. Sample a new vector  $q \sim \text{Dirichlet}(1, \dots, 1)$  of dimension  $k$ .
3. Compute the corrupted probabilities:

$$p_{\text{new}} = \alpha \cdot q + (1 - \alpha) \cdot p.$$

Outcomes with  $p(s'|s, a) = 0$  remain unchanged.

We then run Value Iteration on the corrupted model to compute a policy, which is evaluated on the true environment. This experiment is repeated 10 times for each  $\alpha \in \{0.0, 0.1, \dots, 0.8\}$ , and the mean and standard deviation of the discounted returns are plotted as a function of  $\alpha$ .

### 2.1 Experimental Plot for Task 2

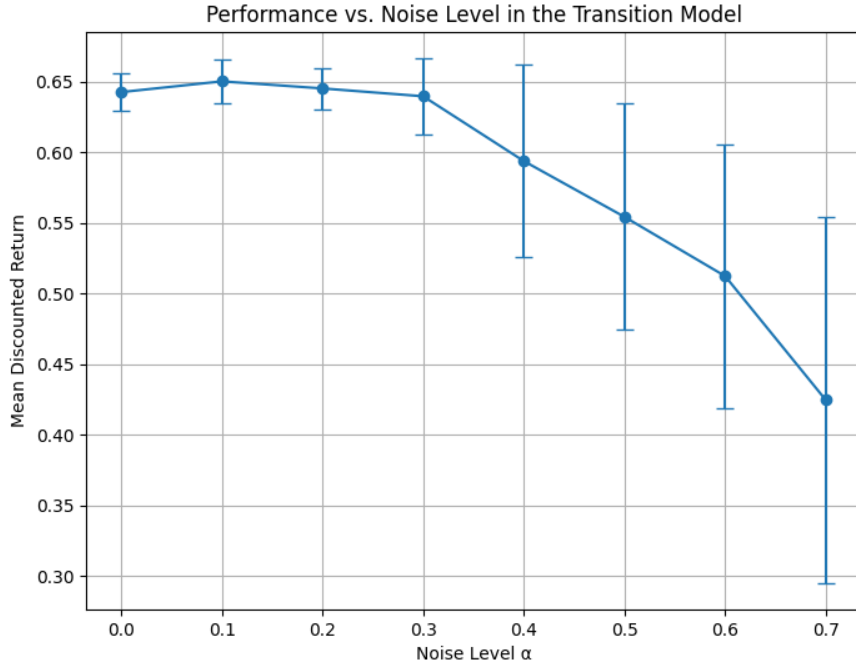


Figure 4: Mean Discounted Return vs. Noise Level  $\alpha$

## 2.2 Observations on Task 2

Our experimental results reveal a trade-off between model accuracy and policy performance. As the noise level  $\alpha$  increases:

- The mean discounted reward drops steeply.
- The variance (and hence the standard deviation) of the reward increases significantly.

These observations indicate that as the transition model becomes less accurate, the resulting policy performance becomes more inconsistent. This suggests that for model-based RL to succeed in this environment, the estimated model must be highly accurate.

### 3 Task 3: Naive Model-Based RL (MBRL)

In this task, we implement a naive model-based RL approach. Data is collected using random actions, and the collected data is used to estimate a model of the environment. Value Iteration (VI) is then run on the estimated model to compute a policy, which is evaluated on the true environment.

#### 3.1 Experimental Plot for Task 3

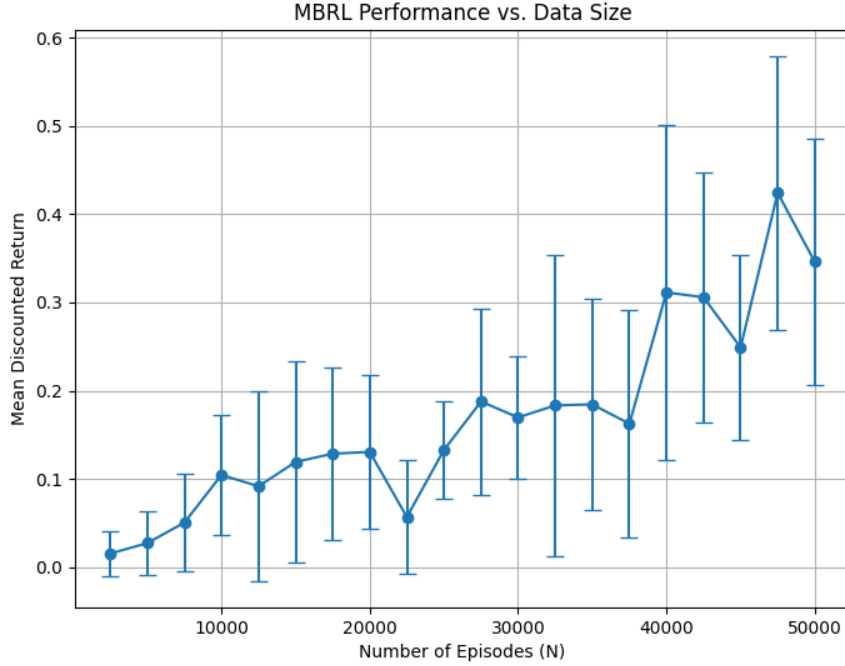


Figure 5: Mean Discounted Return vs. Number of Data Episodes  $N$

#### 3.2 Observations on Task 3

Our results show that as the number of data episodes  $N$  increases, the mean discounted reward improves, indicating a more accurate estimated model and better policy performance. However, the variance (and standard deviation) also increases, meaning that while the average performance is higher, the outcomes become less consistent across experiments.

This suggests that the naive MBRL approach can eventually approach the performance of VI with the correct model, but the higher variance indicates a sensitivity to the particular data samples, implying that additional techniques may be needed to achieve more stable results.

## A Supplementary Terminal Output Results

In this appendix, we show screenshots of the terminal outputs recorded during the execution of the algorithms. These outputs provide additional evidence of the convergence behavior and performance of the algorithms.

### A.1 Value Iteration (VI)

```
--- Value Iteration (VI) ---
Iteration 0: Mean return = 0.0
Iteration 10: Mean return = 0.0
Iteration 20: Mean return = 0.34860819734055715
Iteration 30: Mean return = 0.3386184772173626
Iteration 40: Mean return = 0.3944638706121052
Iteration 50: Mean return = 0.3996334369962461
Iteration 60: Mean return = 0.571425095286304
Iteration 70: Mean return = 0.6295359141881262
Iteration 80: Mean return = 0.5851308125594584
Iteration 90: Mean return = 0.6095084899916863
Iteration 100: Mean return = 0.6444911777288949
Iteration 110: Mean return = 0.6686892408948784
Iteration 120: Mean return = 0.6290049937471672
Iteration 130: Mean return = 0.6064680940177025
Iteration 140: Mean return = 0.6487872944297617
Iteration 150: Mean return = 0.6456605222554723
Iteration 160: Mean return = 0.6372560222864909
Iteration 170: Mean return = 0.67953595942936
Iteration 180: Mean return = 0.6297836514553274
Iteration 190: Mean return = 0.6476314073319802
Iteration 200: Mean return = 0.6658520429242192
Iteration 210: Mean return = 0.6478957514238177
Iteration 220: Mean return = 0.6582796158362875
Iteration 230: Mean return = 0.6361535254046902
Iteration 240: Mean return = 0.6365433958067898
Iteration 250: Mean return = 0.6469033236663583
Converged after 251 iterations.

Final VI Policy (reshaped to 8x8):
[[1 0 0 2 0 0 0 3]
 [3 3 1 2 1 0 3 3]
 [0 2 2 2 1 0 0 2]
 [0 0 2 2 0 0 0 0]
 [0 0 2 2 0 0 0 1]
 [0 0 0 2 3 3 1 3]
 [1 1 1 0 0 0 0 0]
 [1 1 1 1 1 1 1 0]]

Final VI Value Function (reshaped to 8x8):
[[0.5658 0.5681 0.    0.6175 0.6175 0.6159 0.6111 0.6067]
 [0.5681 0.5751 0.5935 0.6218 0.6217 0.6184 0.6106 0.6065]
 [0.    0.3867 0.5879 0.6303 0.629 0.6219 0.    0.3021]
 [0.    0.    0.5438 0.6442 0.6389 0.6222 0.207 0.    ]
 [0.    0.    0.4029 0.6672 0.6476 0.6098 0.    0.1465]
 [0.    0.2967 0.    0.7138 0.6403 0.5631 0.4424 0.2938]
 [0.8943 0.8915 0.8791 0.8374 0.    0.    0.4725 0.    ]
 [0.9009 0.9046 0.912 0.923 0.9375 0.9555 0.9767 0.    ]]
```

Figure 6: Terminal output for Value Iteration (VI)

## A.2 Policy Iteration (PI)

```
--- Policy Iteration (PI) ---
Iteration 0: Mean return = 0.0
Iteration 1: Mean return = 0.0
Iteration 2: Mean return = 0.0
Iteration 3: Mean return = 0.0
Iteration 4: Mean return = 0.0
Iteration 5: Mean return = 0.0
Iteration 6: Mean return = 0.0
Iteration 7: Mean return = 0.18135058268327178
Iteration 8: Mean return = 0.4339775596869663
Iteration 9: Mean return = 0.6386849521974254
Policy converged after 10 iterations

Final PI Policy (reshaped to 8x8):
[[1 0 0 2 0 0 0 3]
 [3 3 1 2 1 0 3 3]
 [0 2 2 2 1 0 0 2]
 [0 0 2 2 0 0 0 0]
 [0 0 2 2 0 0 0 1]
 [0 0 0 2 3 3 1 3]
 [1 1 1 0 0 0 0 0]
 [1 1 1 1 1 1 0 0]]

Final PI Value Function (reshaped to 8x8):
[[0.5581 0.5604 0.    0.6095 0.6095 0.6079 0.6034 0.5993]
 [0.5604 0.5673 0.5856 0.6136 0.6135 0.6103 0.6029 0.599 ]
 [0.    0.3815 0.58    0.6218 0.6205 0.6136 0.    0.2985]
 [0.    0.    0.5363 0.6352 0.6301 0.6138 0.2042 0.    ]
 [0.    0.    0.3972 0.6578 0.6385 0.6016 0.    0.1451]
 [0.    0.2914 0.    0.7034 0.6315 0.5561 0.4383 0.291 ]
 [0.8778 0.8756 0.8646 0.8249 0.    0.    0.4699 0.    ]
 [0.8844 0.8888 0.8976 0.9106 0.9277 0.9487 0.9732 0.    ]]
```

Figure 7: Terminal output for Policy Iteration (PI)

## A.3 Modified Policy Iteration (MPI)

```
--- Modified Policy Iteration (MPI) ---
MPI Iteration 0: Mean return = 0.0000
MPI Iteration 1: Mean return = 0.0000
MPI Iteration 2: Mean return = 0.0000
MPI Iteration 3: Mean return = 0.0000
MPI Iteration 4: Mean return = 0.0000
MPI Iteration 5: Mean return = 0.1792
MPI Iteration 6: Mean return = 0.3198
MPI Iteration 7: Mean return = 0.3662
MPI Iteration 8: Mean return = 0.3053

Final MPI Policy (reshaped to 8x8):
[[1 0 0 2 0 0 0 0]
 [3 2 1 1 1 0 3 3]
 [0 2 2 1 1 0 0 2]
 [0 0 2 2 1 0 0 0]
 [0 0 2 2 1 0 0 1]
 [0 0 0 2 3 3 1 3]
 [1 1 1 0 0 0 0 0]
 [1 2 2 1 1 1 0 0]]

Final MPI Value Function (reshaped to 8x8):
[[0.0046 0.0061 0.    0.0191 0.0197 0.0187 0.013 0.0084]
 [0.0061 0.011 0.0201 0.0265 0.0289 0.0257 0.0137 0.0086]
 [0.    0.0122 0.0299 0.0406 0.0437 0.0391 0.    0.0033]
 [0.    0.    0.0394 0.0594 0.0643 0.0584 0.0181 0.    ]
 [0.    0.    0.0394 0.0884 0.0888 0.0855 0.    0.058 ]
 [0.    0.056 0.    0.1338 0.1095 0.1229 0.1879 0.1195]
 [0.1564 0.1802 0.2137 0.2258 0.    0.    0.3307 0.    ]
 [0.1696 0.2067 0.2694 0.36    0.4838 0.6368 0.8122 0.    ]]
```

Figure 8: Terminal output for Modified Policy Iteration (MPI)



## B README File

```
1 # CSCI B659: Reinforcement Learning - Assignment 1
2 **Author:** LJ Huang
3 **Semester:** Spring 2025
4
5 ## Overview
6 This repository contains the code and report for Assignment 1. The assignment
  covers three main tasks:
7 1. **Task 1:** Implementation and evaluation of Value Iteration (VI), Policy
  Iteration (PI), and Modified Policy Iteration (MPI) on the Frozen Lake
  environment.
8 2. **Task 2:** Exploration of model-based RL by generating corrupted versions of
  the transition model, planning using VI, and evaluating the resulting policies.
9 3. **Task 3:** Implementation of a naive model-based RL approach by collecting
  data with random actions, estimating the environment model, and computing a
  policy using VI.
10
11 ## Directory Structure
12 - **VI_PI_MPI.py**
13   Contains the implementation and experimental driver code for Task 1 (VI, PI, and
   MPI).
14 - **corrupt_version.py**
15   Contains the code for Task 2, including the procedure to corrupt the transition
   model and the associated experiments.
16 - **Naive_MBRL.py**
17   Contains the implementation for Task 3 (naive model-based RL), including data
   collection, model estimation, planning, and evaluation.
18 - **pp1starter.py**
19   The provided startup file for setting up the Frozen Lake environment (including
   reward and transition models, number of states, and actions).
20 - **hw1_report.pdf**
21   The report containing the code printouts, experimental results, plots, and
   discussion of the findings.
22 - **README.md**
23   This file.
24
25 ## Dependencies
26 - Python 3.x
27 - Gymnasium (Install with 'pip install gymnasium')
28 - NumPy (Install with 'pip install numpy')
29 - Matplotlib (Install with 'pip install matplotlib')
30
31 ## Installation and Running
32 1. **Clone or Download the Repository:**
33   Clone the repository or download the zip file and extract its contents.
34
35 2. **Run the Code**
36   python VI_PI_MPI.py (Task 1)
37   python corrupt_version.py (Task 2)
38   python Naive_MBRL.py (Task 3)
```

Listing 1: README.file