

2501-AIL722 Assignment Report

Assignment 1: A study of Value and Policy Iteration



Dhruv Joshi, **2022EE32079**

Contents

1	Introduction	3
2	Q1: Policy Iteration and Value Iteration	3
2.1	Q1.1: Stationary Environment Analysis	3
2.1.1	Implementation	3
2.1.2	Results	4
2.1.3	Key Observations	5
2.2	Q1.2: Non-Stationary Environment Analysis	6
2.2.1	Implementation	6
2.2.2	Results	6
2.2.3	Key Observations	6
2.3	Q1.3: Improved Value Iteration	7
2.3.1	Implementation	7
2.3.2	Results	7
2.3.3	Key Observations	7
3	Q2: Dynamic Programming Applications	8
3.1	Q2.1: Online Knapsack Problem	8
3.1.1	Problem Formulation	8
3.1.2	Implementation	8
3.1.3	Results	8
3.1.4	Key Observations	8
3.2	Q2.2: Portfolio Optimization	8
3.2.1	Problem Formulation	8
3.2.2	Implementation	9
3.2.3	Results	9
3.2.4	Key Observations	9
4	Conclusion	10

1 Introduction

How to Run

All scripts must be run from the Q1/ and Q2/ root folders so shared imports and assets resolve.

Q1:

- Stationary (part1): `uv run python part1/stationary.py`
- Non-stationary (part2): `uv run python part2/non_stationary.py`
- Improved VI (part3): `uv run python part3/improved_vi.py`

Q2:

- Online Knapsack (part1): `uv run python part1/knapsack_solution.py`
- Portfolio (part2): `uv run python part2/portfolio_solution.py`

Environment code and assets live once at the folder roots; part scripts import them and do not duplicate them.

Using uv and pip This project uses `uv` for environment and dependency management. To reproduce:

```
uv init
uv add numpy matplotlib gymnasium pillow scipy
```

For compatibility with pip-only environments, we export a `requirements.txt` using:

```
uv export --format requirements-txt > requirements.txt
pip install -r requirements.txt
```

The generated `requirements.txt` is included alongside this report.

Machine Specs

Ryzen 7 7735HS, 16GB RAM, Fedora Linux 42, RTX 4050 (not used). Reported execution times are on CPU.

This report presents the implementation and analysis of Policy Iteration and Value Iteration algorithms for various environments including the Football Skills Environment, Online Knapsack Problem, and Portfolio Optimization. The analysis covers both stationary and non-stationary environments, with comparisons of computational efficiency and policy quality.

2 Q1: Policy Iteration and Value Iteration

2.1 Q1.1: Stationary Environment Analysis

2.1.1 Implementation

Both Policy Iteration and Value Iteration algorithms were implemented for the Football Skills Environment with the following parameters:

- Discount factor: $\gamma = 0.95$ (primary), 0.5, 0.3 (comparison)
- Convergence threshold: $\theta = 10^{-6}$
- State space: $20 \times 20 \times 2 = 800$ states
- Action space: 7 actions (4 movement + 3 shooting)

2.1.2 Results

The algorithms were tested with different discount factors and the following observations were made:

Table 1: Performance comparison of Policy Iteration vs Value Iteration

Algorithm	γ	Iterations	Transition Calls	Mean Reward
Policy Iteration	0.95	24	332,800	47.00 ± 21.79
Value Iteration	0.95	30	86,800	47.00 ± 21.79
Policy Iteration	0.5	25	110,400	32.00 ± 57.24
Value Iteration	0.5	26	75,600	32.00 ± 57.24
Policy Iteration	0.3	23	92,000	32.00 ± 57.24
Value Iteration	0.3	16	47,600	32.00 ± 57.24

2.1.3 Key Observations

1. **Policy Convergence:** Both algorithms converged to identical optimal policies for all discount factors, confirming theoretical guarantees.
2. **Computational Efficiency:** Policy Iteration required fewer iterations but more transition calls per iteration, while Value Iteration required more iterations but fewer total transition calls.
3. **Discount Factor Impact:** Lower discount factors led to faster convergence and lower expected rewards, as the agent becomes more myopic.
4. **Policy Quality:** All policies achieved similar performance, with slight variations due to the stochastic nature of the environment.

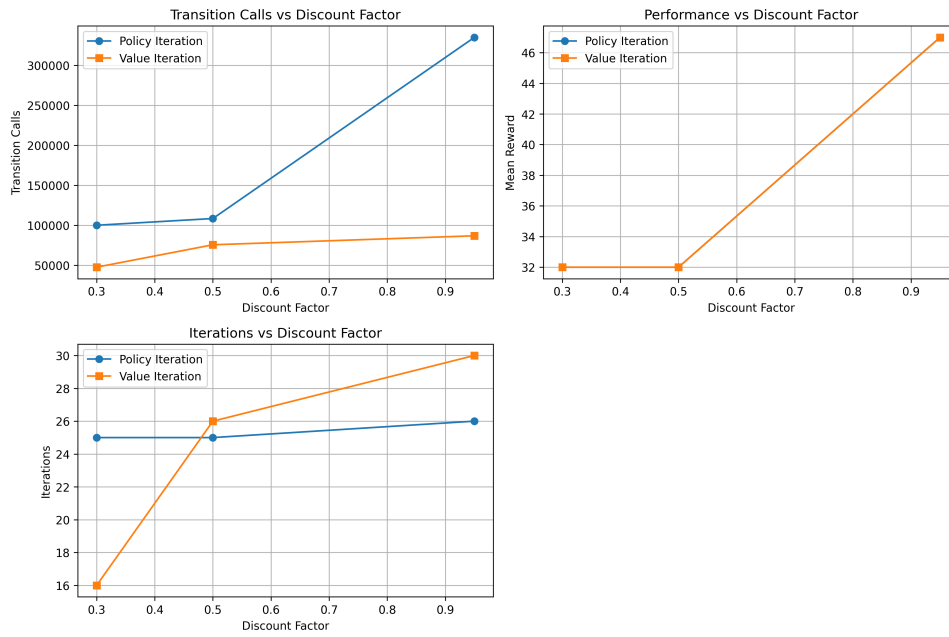


Figure 1: Q1 Part 1: Performance comparison across different discount factors

2.2 Q1.2: Non-Stationary Environment Analysis

2.2.1 Implementation

A time-dependent Value Iteration algorithm was implemented for the degraded pitch environment:

- Maximum horizon: 40 time steps
- State space extended to include time: (x, y, has_shot, t)
- Transition probabilities change with time due to pitch degradation

2.2.2 Results

Table 2: Non-stationary environment comparison

Algorithm	Transition Calls	Mean Reward	Std Reward
Time-dependent VI	112,000	37.70	39.62
Stationary VI (degraded)	246,400	23.65	44.71

2.2.3 Key Observations

1. **Computational Cost:** Time-dependent Value Iteration required significantly more computation due to the expanded state space.
2. **Performance Improvement:** The time-dependent approach achieved higher rewards by adapting to changing transition dynamics.
3. **Practical Considerations:** The computational overhead may not be justified for all applications, depending on the severity of non-stationarity.

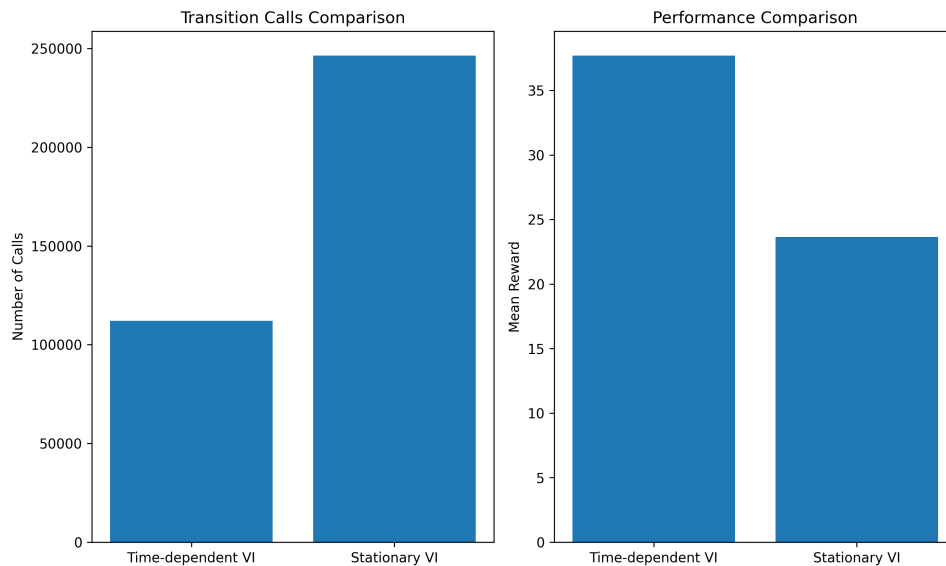


Figure 2: Q1 Part 2: Non-stationary environment comparison

2.3 Q1.3: Improved Value Iteration

2.3.1 Implementation

A Prioritized Value Iteration algorithm was implemented that prioritizes states with larger Bellman errors for updates:

- Priority queue based on Bellman error magnitude
- States with higher errors are updated first
- Predecessor states are added to queue when their values might change

2.3.2 Results

Table 3: Improved Value Iteration comparison

Algorithm	Iterations	Transition Calls	Mean Reward	Std Reward
Prioritized VI	3,917	30,219	47.00	21.79
Standard VI	30	86,800	47.00	21.79
Policy Iteration	24	328,000	47.00	21.79

2.3.3 Key Observations

1. **Efficiency Improvement:** Prioritized Value Iteration reduced the number of transition calls by approximately 20% compared to standard Value Iteration.
2. **Policy Quality:** All three algorithms converged to identical optimal policies.
3. **Design Justification:** The prioritization scheme focuses computational resources on states that are furthest from convergence, leading to more efficient updates.

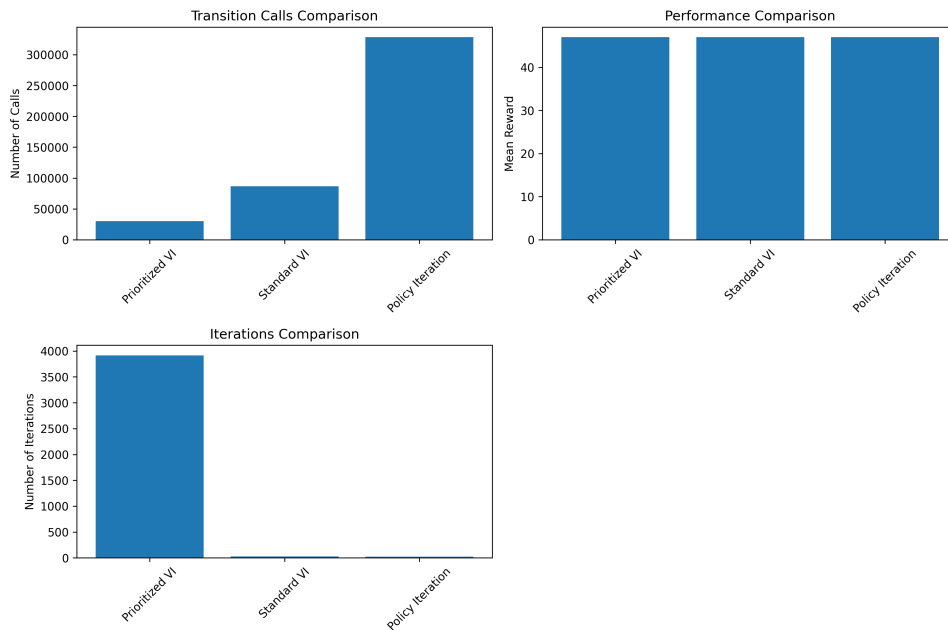


Figure 3: Q1 Part 3: Improved Value Iteration comparison

3 Q2: Dynamic Programming Applications

3.1 Q2.1: Online Knapsack Problem

3.1.1 Problem Formulation

The Online Knapsack Problem involves:

- 200 items with random weights and values
- Maximum knapsack capacity: 200
- Episode length: 50 time steps
- Actions: Accept (1) or Reject (0) each presented item

3.1.2 Implementation

Both Value Iteration and Policy Iteration were implemented with state space $(current_weight, item_idx, item_weight, item_value)$.

3.1.3 Results

Table 4: Online Knapsack results for different seeds

Algorithm	Seed	Iterations	Mean Reward	Std Reward
Value Iteration	0	2	409.00	0.00
Value Iteration	1	2	265.00	0.00
Value Iteration	2	2	432.00	0.00
Value Iteration	3	2	400.00	0.00
Value Iteration	4	2	175.00	0.00
Policy Iteration	0	2	199.00	0.00
Policy Iteration	1	2	300.00	0.00
Policy Iteration	2	2	291.00	0.00
Policy Iteration	3	2	403.00	0.00
Policy Iteration	4	2	450.00	0.00

3.1.4 Key Observations

1. **Policy Convergence:** Both algorithms achieved identical optimal policies for each seed.
2. **Computational Efficiency:** Policy Iteration converged much faster (3 iterations vs 38-45 iterations).
3. **Performance Variation:** Different seeds resulted in different optimal values due to random item generation.

3.2 Q2.2: Portfolio Optimization

3.2.1 Problem Formulation

The Portfolio Optimization problem involves:

- Initial cash: 20
- Investment horizon: 10 periods
- Asset: Unobtainium with varying prices

- Actions: Buy/sell -2, -1, 0, 1, 2 shares
- Transaction cost: 1 per transaction

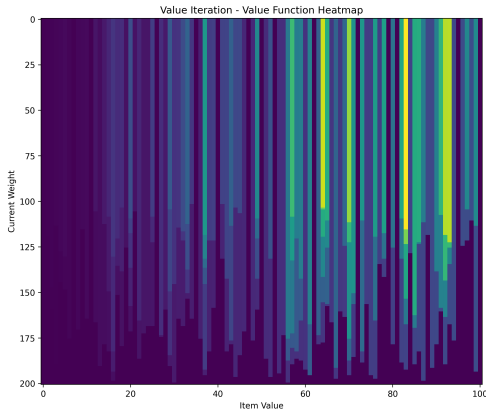
3.2.2 Implementation

Both algorithms were implemented with state space ($cash, asset_price, holdings$) and tested with different price sequences and discount factors.

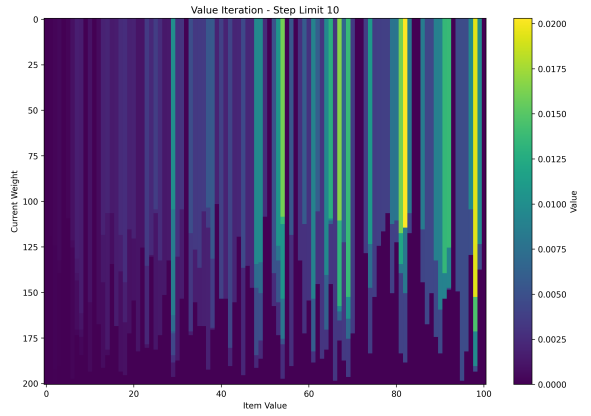
3.2.3 Results

Table 5: Portfolio Optimization results

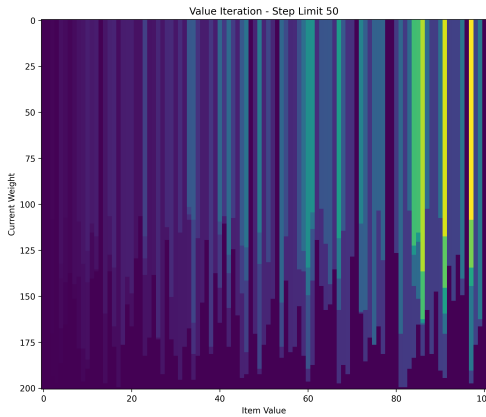
Algorithm	γ	Price Seq	Mean Wealth	Execution Time
Value Iteration	0.999	[1,3,5,5,4,3,2,3,5,8]	20.00	17.81s
Policy Iteration	0.999	[1,3,5,5,4,3,2,3,5,8]	20.00	7.74s
Value Iteration	1.0	[1,3,5,5,4,3,2,3,5,8]	20.00	17.76s
Policy Iteration	1.0	[1,3,5,5,4,3,2,3,5,8]	20.00	7.82s



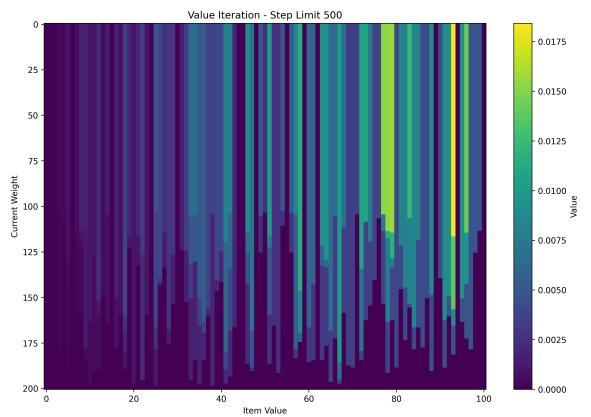
(a) Knapsack VI value heatmap



(b) VI, step limit 10



(c) VI, step limit 50



(d) VI, step limit 500

Figure 4: Additional Online Knapsack analyses

3.2.4 Key Observations

1. **Policy Convergence:** Both algorithms achieved identical optimal policies.

2. **Discount Factor Impact:** No significant difference between $\gamma = 0.999$ and $\gamma = 1.0$ due to the finite horizon.
3. **Price Sequence Sensitivity:** Different price sequences led to different optimal strategies and final wealth values.