

Assignment 2 Report: Escape Through the City

Amirthalingam R
Student ID: 24816

October 12, 2025

1. Implemented Functions (TODOs)

All required methods in `agents.py` were completed. Below are the key implementations.

Value Iteration

```
1 max_q = -np.inf
2 for action in self.env.get_actions():
3     q_value = self._calculate_q_value(state, action)
4     max_q = max(max_q, q_value)
5 new_value_table[state] = max_q
```

Value Iteration - Optimal policy

```
1 for state in self.states:
2     if state == self.env.goal or state in self.env.obstacles:
3         continue
4     best_q = -np.inf
5     for action in self.env.get_actions():
6         q_value = self._calculate_q_value(state, action)
7         if q_value > best_q:
8             best_q = q_value
9             current_policy[state] = action
```

Policy Evaluation

```
1 new_value_table[state] = self._calculate_q_value(state, self.policy[state])
```

Policy Improvement

```
1 for action in self.env.get_actions():
2     if self._calculate_q_value(state, action) > self._calculate_q_value(state, best_action):
3         best_action = action
4 self.policy[state] = best_action
```

2. Results

Stage 1: Sunny Day (No Obstacles, No Rain)

Convergence behaviour:

Discount factor	Iterations	Reward	Converged ?	Correct Policy ?
0.9	19	-11	Yes	Yes
0.7	19	-11	Yes	Yes
0.5	19	-11	Yes	Yes
0.3	13	-11	Yes	Yes
0.1	7	-100	Yes	No

Table 1: Convergence behavior of Value Iteration

Discount factor	Iterations	Reward	Converged ?	Correct Policy ?
0.9	19	-11	Yes	Yes
0.7	19	-11	Yes	Yes
0.5	19	-11	Yes	Yes
0.3	19	-11	Yes	Yes
0.1	17	-11	Yes	Yes

Table 2: Convergence behavior of Policy Iteration

Observations:

- Both algorithms converged to the same optimal policy in all cases except when the discount factor was 0.1.
- Value Iteration did not learn the correct policy when the γ value is small. This is because the agent becomes myopic and does not provide enough weightage to long term rewards.
- Policy Iteration always converged and learned the optimal policy.
- Lower γ gives faster convergence. But, can sometimes result in learning sub optimal policy.

Learned Policies

[illegible]Figure 1: Value Iteration, $\gamma=0.9$ [illegible]Figure 2: Policy Iteration, $\gamma=0.9$ [illegible]Figure 3: Value Iteration, $\gamma=0.7$ [illegible]Figure 4: Policy Iteration, $\gamma=0.7$

Agent's Learned Policy									
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	\$	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	→	→	→	→	→	→	→	→	G

Figure 5: Value Iteration, $\gamma=0.5$

Agent's Learned Policy									
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	\$	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	→	→	→	→	→	→	→	→	G

Figure 6: Policy Iteration, $\gamma=0.5$

Agent's Learned Policy									
↑	↑	↑	↑	↑	↓	↓	↓	↓	↓
↑	↑	↑	↑	↓	↓	\$	↓	↓	↓
↑	↑	↑	↓	↓	↓	↓	↓	↓	↓
↑	↑	↓	↓	↓	↓	↓	↓	↓	↓
↑	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	→	→	→	→	→	→	→	→	G

Figure 7: Value Iteration, $\gamma=0.3$

Agent's Learned Policy									
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	\$	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	→	→	→	→	→	→	→	→	G

Figure 8: Policy Iteration, $\gamma=0.3$

Agent's Learned Policy									
↑	↑	↑	↑	↑	↑	↑	↑	↑	↑
↑	↑	↑	↑	↑	↑	\$	↑	↑	↑
↑	↑	↑	↑	↑	↑	↑	↑	↑	↓
↑	↑	↑	↑	↑	↑	↑	↑	↓	↓
↑	↑	↑	↑	↑	↑	↑	↓	↓	↓
↑	↑	↑	↑	↑	↑	↓	↓	↓	↓
↑	↑	↑	↑	↑	↓	↓	↓	↓	↓
↑	↑	↑	↑	↓	↓	↓	↓	↓	↓
↑	↑	↑	↓	↓	↓	↓	↓	↓	↓
↑	↑	→	→	→	→	→	→	→	G

Figure 9: Value Iteration, $\gamma=0.1$

Agent's Learned Policy									
↑	↑	↓	↓	↓	↓	↓	↓	↓	↓
↑	↓	↓	↓	↓	↓	\$	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	→	→	→	→	→	→	→	→	G

Figure 10: Policy Iteration, $\gamma=0.1$

Stage 2: Obstacles Added

Convergence behaviour:

Discount factor	Iterations	Reward	Converged ?	Correct Policy ?
0.9	19	-11	Yes	Yes
0.7	19	-11	Yes	Yes
0.5	19	-11	Yes	Yes
0.3	13	-11	Yes	Yes
0.1	7	-100	Yes	No

Table 3: Convergence behavior of Value Iteration

Discount factor	Iterations	Reward	Converged ?	Correct Policy ?
0.9	19	-11	Yes	Yes
0.7	19	-11	Yes	Yes
0.5	19	-11	Yes	Yes
0.3	19	-11	Yes	Yes
0.1	17	-11	Yes	Yes

Table 4: Convergence behavior of Policy Iteration

Observations:

- Both algorithms successfully avoided obstacle cells. Did not note any difference in convergence speed when compared to the previous stage.
- Value iteration with $\gamma=0.1$ did not learn the optimal policy.
- In other cases, the agent learned the shortest path while still avoiding the obstacles.

Learned Policies

Agent's Learned Policy

↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	\$	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	↓	↓	↓	↓	↓	→	↓	↓	↓
█	↓	↓	→	↓	↓	█	↓	↓	↓
↓	→	↓	█	↓	↓	↓	↓	→	↓
↓	█	↓	↓	↓	↓	↓	↓	█	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	→	→	→	→	→	→	→	→	G

Figure 11: Value Iteration, $\gamma=0.9$

Agent's Learned Policy

↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	\$	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	↓	↓	↓	↓	↓	→	↓	↓	↓
█	↓	↓	→	↓	↓	█	↓	↓	↓
↓	→	↓	█	↓	↓	↓	↓	→	↓
↓	█	↓	↓	↓	↓	↓	↓	█	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	→	→	→	→	→	→	→	→	G

Figure 12: Policy Iteration, $\gamma=0.9$

Agent's Learned Policy

↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	\$	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	↓	↓	↓	↓	↓	→	↓	↓	↓
█	↓	↓	→	↓	↓	█	↓	↓	↓
↓	→	↓	█	↓	↓	↓	↓	→	↓
↓	█	↓	↓	↓	↓	↓	↓	█	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	→	→	→	→	→	→	→	→	G

Figure 13: Value Iteration, $\gamma=0.7$

Agent's Learned Policy

↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	\$	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	↓	↓	↓	↓	↓	→	↓	↓	↓
█	↓	↓	→	↓	↓	█	↓	↓	↓
↓	→	↓	█	↓	↓	↓	↓	→	↓
↓	█	↓	↓	↓	↓	↓	↓	█	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	→	→	→	→	→	→	→	→	G

Figure 14: Policy Iteration, $\gamma=0.7$

Agent's Learned Policy									
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	\$	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	↓	↓	↓	↓	↓	→	↓	↓	↓
	↓	↓	→	↓	↓		↓	↓	↓
↓	→	↓		↓	↓	↓	↓	→	↓
↓		↓	↓	↓	↓	↓	↓		↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	→	→	→	→	→	→	→	→	G

Figure 15: Value Iteration, $\gamma=0.5$

Agent's Learned Policy									
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	\$	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	↓	↓	↓	↓	↓	→	↓	↓	↓
	↓	↓	→	↓	↓		↓	↓	↓
↓	→	↓		↓	↓	↓	↓	→	↓
↓		↓	↓	↓	↓	↓	↓		↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	→	→	→	→	→	→	→	→	G

Figure 16: Policy Iteration, $\gamma=0.5$

Agent's Learned Policy									
↑	↑	↑	↑	↑	↓	↓	↓	↓	↓
↑	↑	↑	↑	↓	↓	\$	↓	↓	↓
↑	↑	↑	↓	↓	↓	↓	↓	↓	↓
↑	↑	↓	↓	↓	↓	→	↓	↓	↓
	↓	↓	→	↓	↓		↓	↓	↓
↓	→	↓		↓	↓	↓	↓	→	↓
↓		↓	↓	↓	↓	↓	↓		↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	→	→	→	→	→	→	→	→	G

Figure 17: Value Iteration, $\gamma=0.3$

Agent's Learned Policy									
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	\$	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	↓	↓	↓	↓	↓	→	↓	↓	↓
	↓	↓	→	↓	↓		↓	↓	↓
↓	→	↓		↓	↓	↓	↓	→	↓
↓		↓	↓	↓	↓	↓	↓		↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	→	→	→	→	→	→	→	→	G

Figure 18: Policy Iteration, $\gamma=0.3$

Agent's Learned Policy

↑	↑	↑	↑	↑	↑	↑	↑	↑	↑
↑	↑	↑	↑	↑	↑	\$	↑	↑	↑
↑	↑	↑	↑	↑	↑	↑	↑	↑	↓
↑	↑	↑	↑	↑	↑	↑	↑	↓	↓
█	↑	↑	↑	↑	↑	█	↓	↓	↓
↓	↑	↑	█	↑	↑	↓	↓	→	↓
↑	█	↑	↓	↑	↓	↓	↓	█	↓
↑	↓	↑	↑	↓	↓	↓	↓	↓	↓
↑	↑	↑	↓	↓	↓	↓	↓	↓	↓
↑	↑	→	→	→	→	→	→	→	G

Figure 19: Value Iteration, $\gamma=0.1$

Agent's Learned Policy

↑	↑	↓	↓	↓	↓	↓	↓	↓	↓
↑	↓	↓	↓	↓	↓	\$	↓	↓	↓
↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
→	↓	↓	↓	↓	↓	→	↓	↓	↓
█	↓	↓	→	↓	↓	█	↓	↓	↓
→	→	↓	█	↓	↓	↓	↓	→	↓
↓	█	↓	↓	↓	↓	↓	↓	█	↓
→	→	↓	↓	↓	↓	↓	↓	↓	↓
→	→	↓	↓	↓	↓	↓	↓	↓	↓
→	→	→	→	→	→	→	→	→	G

Figure 20: Policy Iteration, $\gamma=0.1$

Stage 3: Rainy Day (Stochastic Slips)

Convergence behaviour:

Slip prob	Theta	Discount factor	Iterations	Reward	Converged ?
0.0001	1e-6	0.8	21	-11	Yes
0.001	1e-6	0.8	92	-11	Yes
0.01	1e-3	0.8	229	-11	Yes

Table 5: Convergence behavior of Value Iteration

Slip prob	Theta	Discount factor	Iterations	Reward	Converged ?
0.0001	1e-6	0.8	10	-11	Yes
0.001	1e-6	0.8	9	-11	Yes
0.01	1e-3	0.8	235	-11	Yes

Table 6: Convergence behavior of Policy Iteration

Observations:

- As slip probability increased, convergence slowed down noticeably.
- Value Iteration remained more stable, while Policy Iteration oscillated.
- For slip = 0.01, both Value Iteration and Policy Iteration converged with higher θ tolerance only (1e-3).

Learned Policies

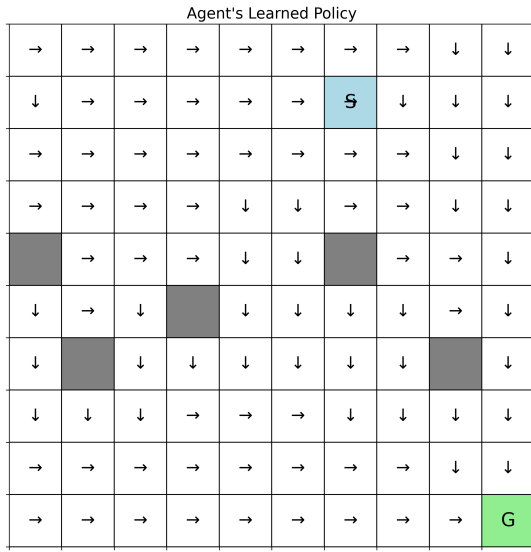


Figure 21: Value Iteration, $\text{slip}=0.0001$

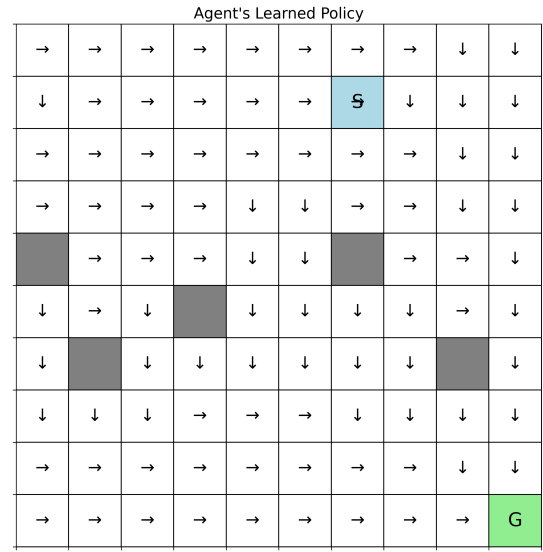


Figure 22: Policy Iteration, $\text{slip}=0.0001$

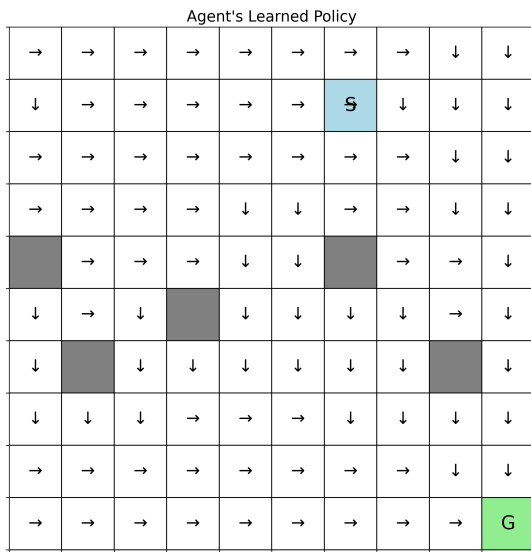


Figure 23: Value Iteration, $\text{slip}=0.001$

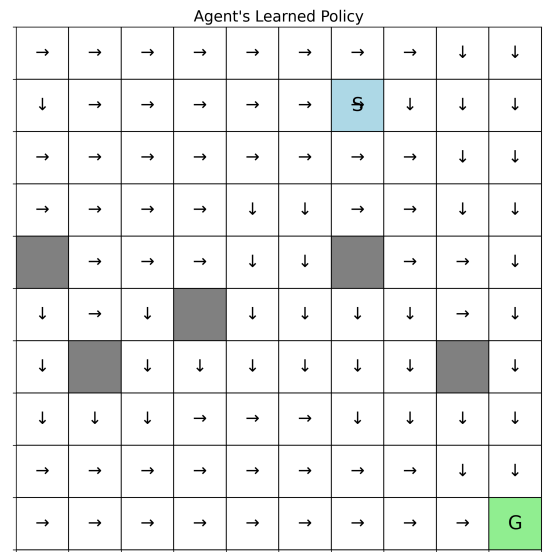


Figure 24: Policy Iteration, $\text{slip}=0.001$

Agent's Learned Policy									
→	→	→	→	→	→	→	→	↓	↓
→	→	→	→	→	→	S	↓	↓	↓
→	→	→	→	→	→	→	→	↓	↓
→	→	→	→	↓	↓	→	→	↓	↓
■	→	→	→	↓	↓	■	→	→	↓
↓	→	↓	■	↓	↓	↓	↓	→	↓
↓	■	↓	↓	↓	↓	↓	↓	■	↓
↓	↓	↓	→	→	→	↓	↓	↓	↓
→	→	→	→	→	→	→	→	↓	↓
→	→	→	→	→	→	→	→	→	G

Figure 25: Value Iteration, slip=0.01

Agent's Learned Policy									
→	→	→	→	→	→	→	→	↓	↓
→	→	→	→	→	→	S	↓	↓	↓
→	→	→	→	→	→	→	→	↓	↓
→	→	→	→	↓	↓	→	→	↓	↓
■	→	→	→	↓	↓	■	→	→	↓
↓	→	↓	■	↓	↓	↓	↓	→	↓
↓	■	↓	↓	↓	↓	↓	↓	■	↓
↓	↓	↓	→	→	→	↓	↓	↓	↓
→	→	→	→	→	→	→	→	↓	↓
→	→	→	→	→	→	→	→	→	G

Figure 26: Policy Iteration, slip=0.01

3. Insights

- **Discount factor (γ):** Lower γ values led to faster but greedier convergence; higher γ yielded more optimal long-term paths.
- **Obstacles:** Forced exploration of alternative routes; both algorithms handled penalties well.
- **Stochasticity:** Introduced instability; tuning θ improved convergence for noisy environments.
- **Performance:** Policy Iteration converged faster in deterministic cases, while Value Iteration (with higher γ) was more robust under uncertainty.

4. Summary

- Both algorithms produced optimal paths in all stages except value iteration with lower γ .
- Value Iteration (with higher γ) offered better stability with noise and uncertainty.
- Policy Iteration was efficient when transitions were deterministic.
- Increasing complexity (obstacles + rain) demonstrated how dynamic programming scales to uncertain environments.