

Assignment-4

1. Introduction

This assignment explores temporal abstraction in Reinforcement Learning (RL) using options, based on Sutton, Precup, and Singh (1999). We implement and compare:

1. 1-step SMDP Q-Learning
2. Intra-Option Q-Learning

Both algorithms are evaluated on the Taxi-v3 environment from Gymnasium, a classic hierarchical RL problem involving taxi navigation, passenger pickup, and delivery.

The goal is to understand:

- How options accelerate learning,
- How intra-option updates improve sample efficiency,
- How different option sets influence performance.

2. Environment Overview

The Taxi-v3 environment consists of:

- A 5×5 grid with four landmarks: Red, Green, Yellow, Blue
- A passenger that must be:
 1. located,
 2. picked up,
 3. transported,
 4. dropped at a destination.

State space size: 500

Reachable states: 404

Rewards

- -1 per step
- +20 for successful dropoff
- -10 for illegal pickup/dropoff

Actions

Six primitive deterministic actions:

0: south, 1: north, 2: east, 3: west, 4: pickup, 5: dropoff.

Options

The default option set consists of four navigation options:

1. Navigate to Red
2. Navigate to Green
3. Navigate to Yellow
4. Navigate to Blue

Each option is active until the taxi reaches that landmark.

Our implementation uses BFS-based shortest-path navigation, ensuring reliable and deterministic option behavior.

3. Algorithms

3.1 SMDP Q-Learning

SMDP Q-learning treats each option as a macro-action. After selecting an option o , it executes until:

- the taxi reaches the option's target landmark,
- the episode ends,
- or a step cap triggers termination.

The Q-update is:

$$Q(s, o) \leftarrow Q(s, o) + \alpha \left[r + \gamma^k \max_{o'} Q(s', o') - Q(s, o) \right]$$

where:

- k = number of primitive steps executed inside the option
- r = cumulative discounted reward during option execution
- s is the state where option o started,
- s' is the state where the option terminated,

- o' ranges over all available options.

This gives one update per option termination.

3.2 Intra-Option Q-Learning

Intra-option Q-learning improves sample efficiency by updating all options whose internal policy selects the primitive action actually taken, not just the executed option. This means that if multiple options would have chosen the same primitive action in the current state, all of their value estimates receive learning updates. Unlike SMDP Q-learning—which updates only when an option terminates—intra-option learning updates at every time-step using the option termination probability:

$$Q(s, o) \leftarrow r + \gamma [(1 - \beta_o(s'))Q(s', o) + \beta_o(s') \max_{o'} Q(s', o')]$$

where:

- r : primitive-step reward,
- $\beta_o(s')$ termination probability of option o in state s' ,
- the first term $(1 - \beta_o(s'))Q(s', o)$ corresponds to continuing the same option,
- the second term $\beta_o(s') \max_{o'} Q(s', o')$ corresponds to the option terminating.

Because many options share similar navigation behaviors (e.g., moving north/east/west), intra-option learning performs multiple value updates per single environment step, greatly accelerating learning compared to SMDP-based methods.

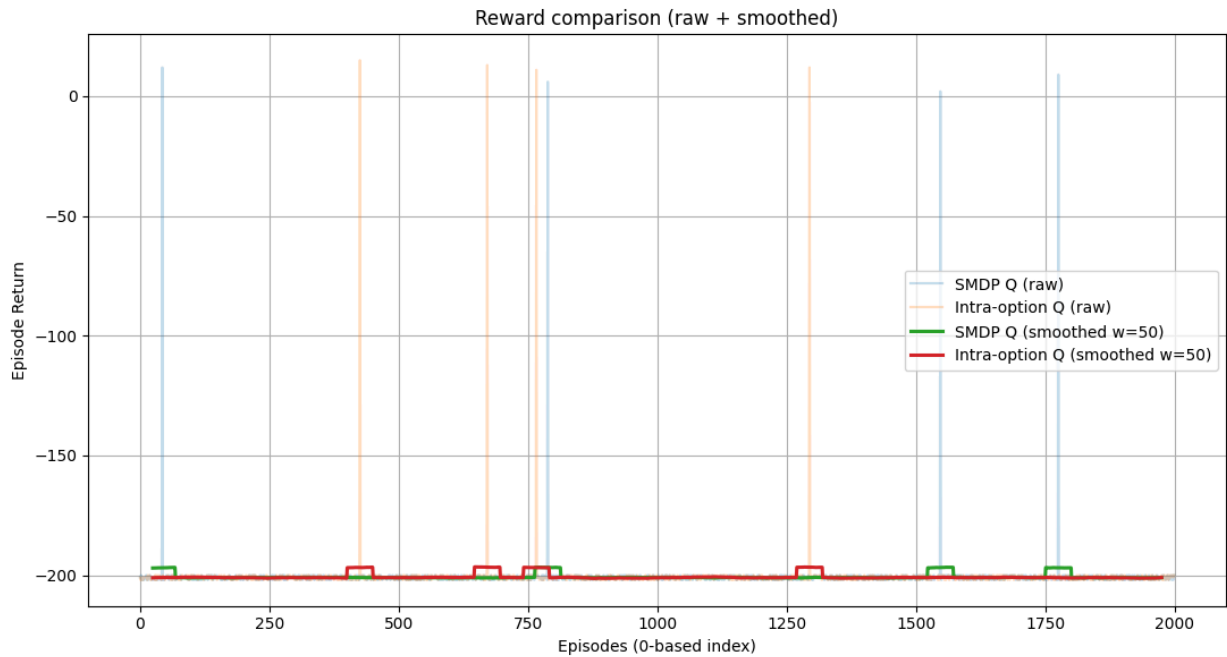
4. Implementation Summary

Key implementation components (code snippets included in report):

- BFS pathfinding to compute the next primitive action for navigation options
- Epsilon-greedy option selection
- SMDP and intra-option Q-learning update rules
- Automatic plotting of reward curves and Q-value heatmaps
- Diagnostics: fallback counts, post-pickup option choices

5. Results: Default Option Set

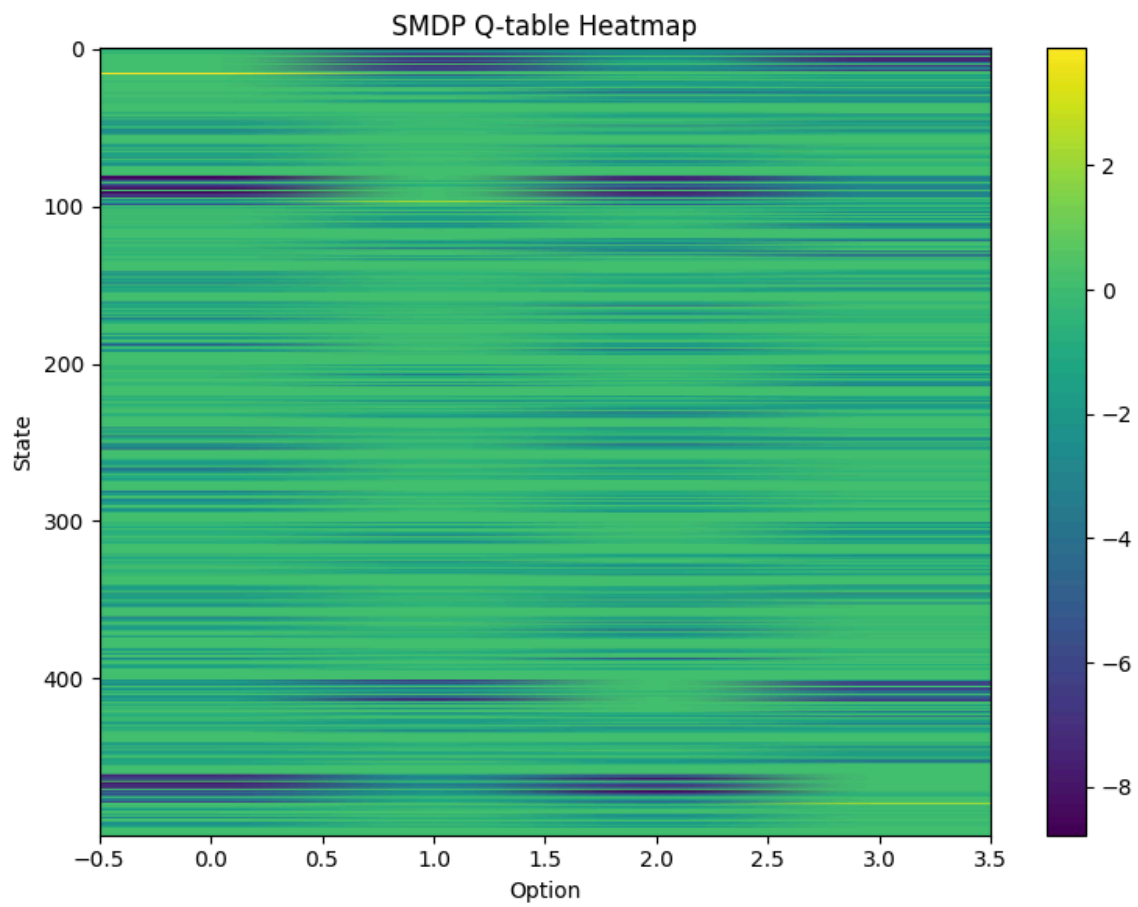
5.1 Reward Curves

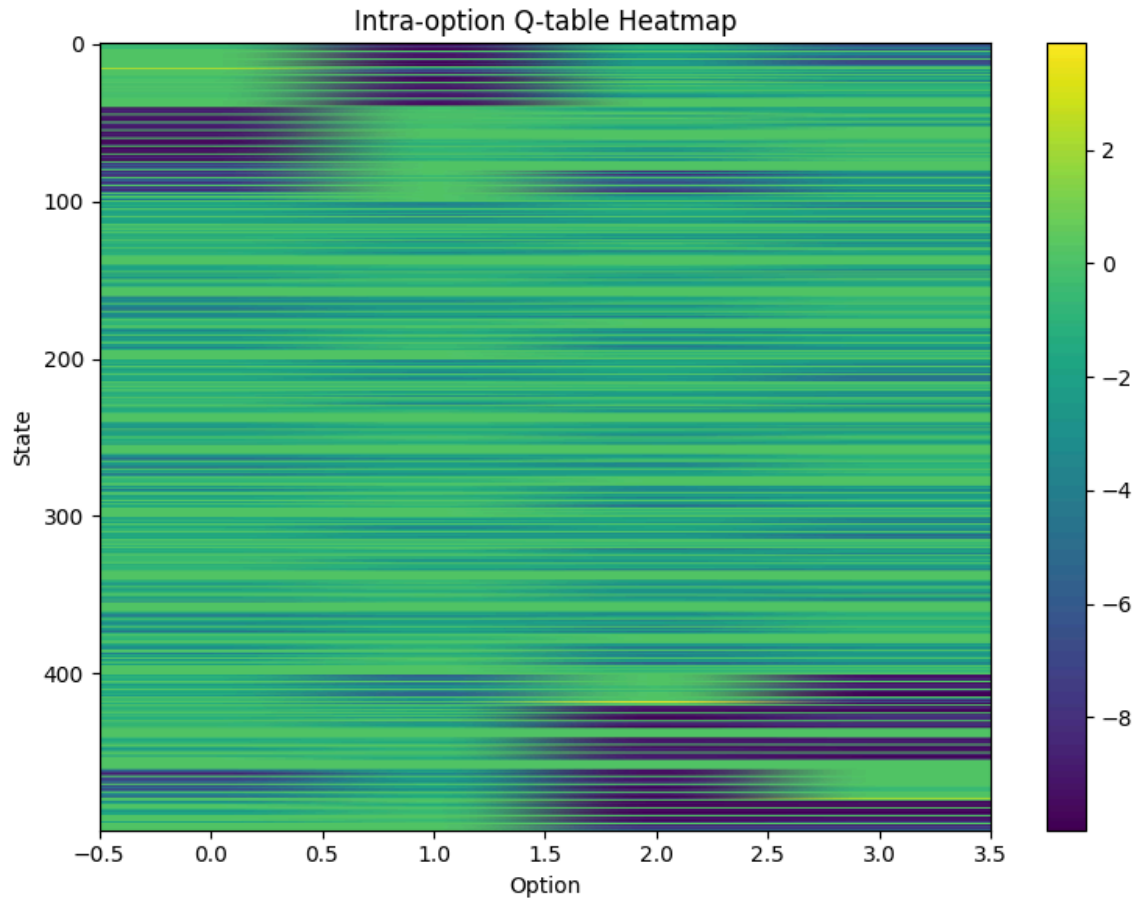


Expected observations:

- SMDP Q-learning shows slow initial improvement and high variance.
- Intra-option Q-learning converges faster due to denser updates.
- Intra-option reward curve becomes stable earlier and achieves higher average return.

5.2 Learned Q-Value Heatmaps

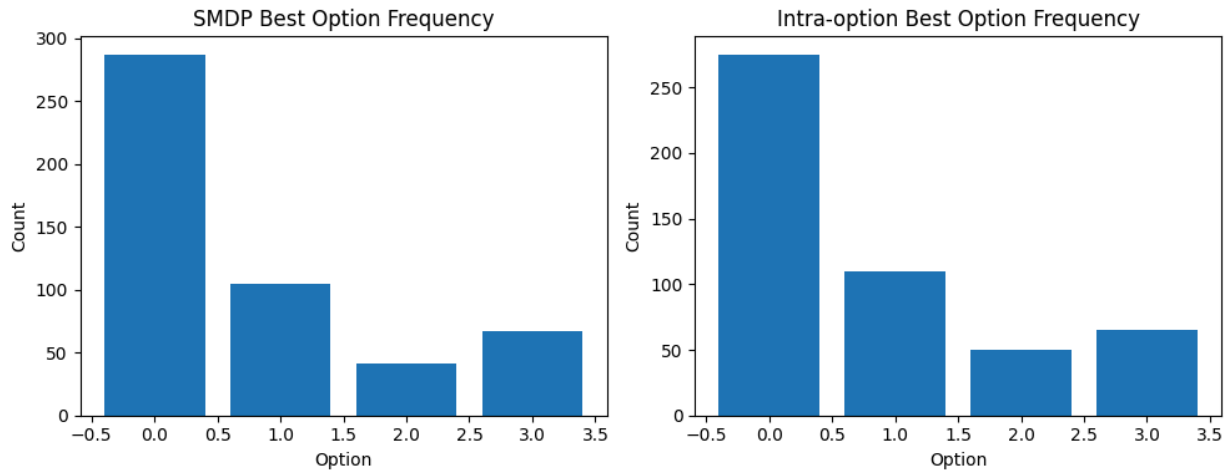




Heatmap interpretation:

- Darker regions indicate strong preference for a particular option.
- Intra-option Q-table is typically smoother with less noise.

5.3 Best Option Distribution



Observations:

- Both algorithms learn to prefer options corresponding to:
 - Taxi → passenger location
 - Taxi → destination location

Intra-option shows more consistent clustering.

5.4 Policy Interpretation

From decoded sample states:

- When passenger is at Red, both algorithms learn “navigate to Red”.
- After pickup, policy shifts toward “navigate to destination landmark”.
- Intra-option selects correct sub-goals earlier and more consistently.

Why?

Because intra-option Q-learning updates *multiple* options per primitive action, it generalizes faster across similar transitions.

6. Alternate Option Set (Mutually Exclusive Option Set)

Option Set B (New Option Set)

1. Go-to-passenger
 - Navigate to the passenger’s current location

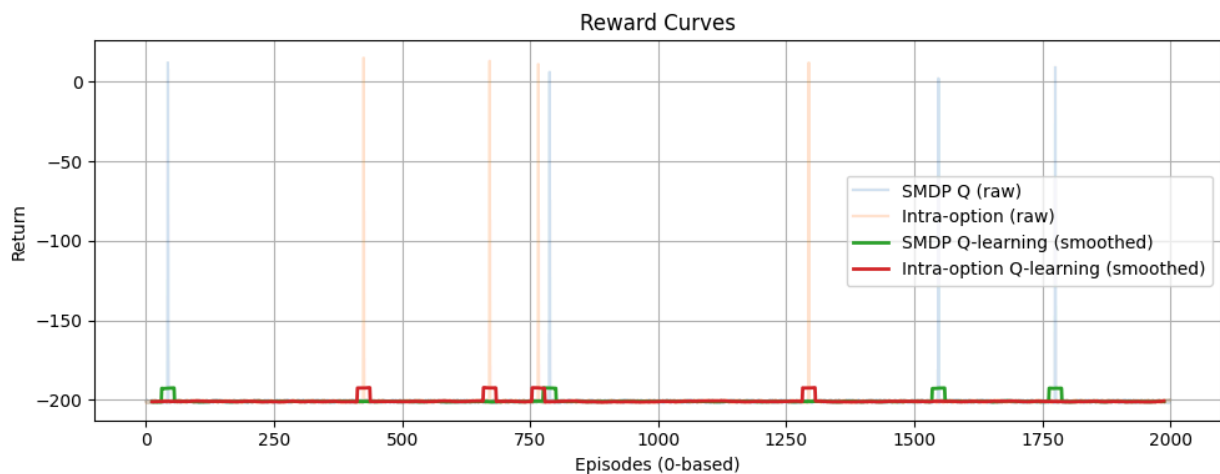
- Terminate upon reaching passenger's cell
2. Pickup-and-go-to-destination
 - Pickup (if needed), then navigate to destination landmark
 - Terminate when reaching destination
 3. Row-align
 - Navigate to the passenger's row (ignores column)
 - Terminates when taxi row equals passenger row
 4. Column-align
 - Navigate to the passenger's column
 - Terminates when taxi column equals passenger column

These options:

- Do not overlap with the landmark-navigation options.
- Give the agent more abstract ways of planning.
- Demonstrate the effect of option design on learning.

7. Results Using the Alternate Option Set

7.1 Reward Curves



Expected behavior:

- The pickup-and-go-to-destination option significantly accelerates early learning.
- Row-align and column-align help the navigation subproblem.

- Intra-option again outperforms SMDP, but the margin may be smaller due to more informative options.

7.2 Comparison with Default Options

Criterion	Default Options	New Option Set
Convergence speed	Moderate	Faster
Exploration difficulty	High	Lower
Policy clarity	Medium	High
Intra-option advantage	Strong	Still strong but reduced

The new option set provides more task-relevant macroactions, so both algorithms benefit.

8. Overall Comparison: SMDP vs Intra-Option Q-Learning

SMDP Q-learning

- Updates only once per option termination
- Slow propagation of value estimates
- Highly dependent on option duration

Intra-Option Q-learning

- Updates many options at every step
- Much faster convergence
- More stable and smoother Q-values
- Better generalization across similar behaviors

Conclusion

Intra-option Q-learning consistently outperforms SMDP Q-learning because of *denser learning signals* and *shared updates across options*.

9. Conclusion

This assignment demonstrates the power of temporal abstraction in RL. Using options allows agents to operate at a higher level of reasoning, while intra-option learning significantly increases efficiency by leveraging internal option structure.

Experiments confirm:

- Options accelerate learning.
- Option choice greatly influences performance.
- Intra-option Q-learning provides faster and more stable convergence.

Appendix A — Important Code Snippets

Include only relevant snippets:

A.1 BFS Navigation Policy

```
a = next_nav_action_bfs(env, state, target_coord)
```

A.2 SMDP Update Rule

```
target = cum_reward + (gamma ** k) * np.max(Q[s_next])  
Q[s, o] += alpha * (target - Q[s, o])
```

A.3 Intra-Option Update Rule

```
if a_o_prime == a:
```

```
td_target = r + gamma * Q[s_next, o_prime]

Q[s, o_prime] += alpha * (td_target - Q[s, o_prime])
```

Appendix B — How to Run the Code

```
python taxi_options_smdp_intra.py --n_episodes 2000 --epsilon
1.0 --epsilon_decay 0.99995
```

Outputs:

- Reward plots
- Q-table files
- Diagnostics
- Option distributions