

Tabular Q-Learning

Total Points: 45

Part 1: Tabular Q-Learning

Task 1: TaxiAgent Implementation (15 points)

Full marks for this task: 15 points

Item	Point s	Criteria
Q-table initialization	3	Q-table is initialized as a 2D NumPy array with shape <code>(env.observation_space.n, env.action_space.n)</code> using the provided <code>q_init</code> parameter
Epsilon greedy action selection	3	<code>action_select()</code> method returns <code>env.action_space.sample()</code> with probability <code>epsilon</code> , otherwise returns <code>np.argmax(self.Q[s])</code>
Q-learning update - target calculation	3	<code>update()</code> method correctly computes target as <code>r</code> when <code>terminated</code> is True, and as <code>r + self.discount * np.max(self.Q[s_new])</code> when <code>terminated</code> is False
Q-learning update - Q-table update	3	<code>update()</code> method correctly applies the update rule: <code>self.Q[s_old, a] += self.lr * (target - self.Q[s_old, a])</code>
Epsilon decay	3	<code>update()</code> method decreases epsilon by multiplying by <code>epsilon_decay</code> while ensuring epsilon does not go below <code>epsilon_final</code> (using <code>max()</code> or equivalent logic)

Task 2: Training Implementation and Results (15 points)

Full marks for this task: 15 points

Item	Point s	Criteria
Training loop structure	3	<code>train()</code> function includes: environment reset, action selection using <code>agent.action_select()</code> , environment step, and Q-learning update using <code>agent.update()</code> for each timestep
Episode termination handling	3	Training loop correctly handles episode termination by checking both <code>terminated</code> and <code>truncated</code> flags and resetting the environment when either is True
Training execution	3	Code successfully runs training for the specified number of episodes without errors
Learning curve visualization	3	Code generates both plots (episode returns and episode lengths) using the provided <code>visualize_learning_curves()</code> function
Convergence achievement	3	Final average return over the last 100 episodes is between -5 and +15 (indicating reasonable learning occurred)

Task 3: Reflection and Analysis Questions (15 points)

Full marks for this task: 15 points

Question 1 (5 points)

Item	Point s	Criteria
Reason 1	2	Identifies Q-table initialization issue
Reason 2	2	Mentions high epsilon exploration
Reason 3	1	References invalid actions penalty

Question 2 (5 points)

Item	Point s	Criteria
States convergence range	2	Provides a specific episode length convergence value or range
Explains optimal path reasoning	2	Explains that optimal pickup and delivery requires a certain minimum number of moves
Connects episode length to returns	1	Notes the relationship between shorter episodes and higher returns

Question 3 (5 points)

Item	Point s	Criteria
Explains Q-value gradient	2	Describes how explored (s,a) pairs get updated
References action selection mechanism	2	Explains that <code>argmax</code> during exploitation will prefer unexplored actions
Mentions update rule specifics	1	References the Q-learning update equation and how the -1 reward propagates through the Q-table via the update rule

Optional Visualization Bonus

This is optional bonus credit for implementing the qualitative policy visualization

Item	Point s	Criteria
Environment setup	1	Creates a new Taxi environment with render mode (either "human" or "ansi") and sets agent epsilon to 0 or a low value for evaluation
Evaluation loop structure	1	Implements a loop that runs multiple episodes (typically 3-5) with the trained agent, including environment reset, action selection, and episode termination handling

Results display	1	Displays or prints information about each episode such as the environment state (using <code>env.render()</code>), total reward, or steps taken
-----------------	---	--

Deep Reinforcement Learning Assignment Rubric

Total Points: 55

Part 2: Deep Reinforcement Learning

Task 1: Neural Network Architecture (12 points)

Full marks for this task: 12 points

Item	Point s	Criteria
Input layer dimension	3	Input layer (<code>fc1</code> or equivalent) accepts 8-dimensional input (first <code>nn.Linear</code> has input size 8)
Hidden layer structure	3	At least 2 fully connected hidden layers are defined with at least 64 neurons each
ReLU activation	3	ReLU (or equivalent non-linear) activation function is applied between layers (not on output layer)
Output layer dimension	3	Output layer produces 4-dimensional output (final <code>nn.Linear</code> has output size 4) with no activation function applied

Task 2: LunarLanderAgent Implementation (21 points)

Full marks for this task: 21 points

Item	Point s	Criteria
------	---------	----------

Two network initialization	3	Both <code>q_network</code> and <code>target_network</code> are instantiated using the <code>NeuralNet</code> class and moved to the correct device
Target network parameter copying	3	Target network parameters are initialized by copying from the main Q-network using <code>load_state_dict()</code> or equivalent
Optimizer and loss initialization	3	An optimizer (e.g., Adam) is initialized with <code>q_network.parameters()</code> and a loss function (e.g., <code>MSELoss</code>) is defined
Replay buffer implementation	3	Replay buffer is implemented using <code>deque</code> with specified <code> maxlen</code> , and experiences (state, action, reward, next_state, terminated) are appended in <code>update()</code> method
Epsilon-greedy action selection	3	<code>action_select()</code> returns random action with probability epsilon, otherwise converts state to tensor, passes through <code>q_network</code> , and returns action with maximum Q-value using <code>argmax()</code>
Q-value computation	3	Current Q-values are computed by passing states through <code>q_network</code> and selecting values for taken actions using <code>gather()</code> or equivalent indexing
Target Q-value computation	3	Target Q-values are computed using <code>target_network</code> with <code>torch.no_grad()</code> , using the formula: <code>rewards + (1 - terminated) * discount * max(target_network(next_states))</code>

Task 3: Training and Results (12 points)

Full marks for this task: 12 points

Item	Point s	Criteria
Training loop execution	3	Training code runs for at least 500 episodes without errors and includes action selection, environment step, and agent update
Terminated vs truncated handling	3	Only <code>terminated</code> flag (not <code>truncated</code>) is passed to <code>agent.update()</code> as the done signal

Environment interaction structure	3	Training loop includes: environment reset at start of each episode, action selection using <code>agent.action_select()</code> , environment step with <code>env.step()</code> , and agent update using <code>agent.update()</code> for each timestep
Learning outcome	3	The average/expected episode return converges to greater than 150

Task 4: Reflection and Analysis Questions (10 points)

Full marks for this task: 10 points

Question 1 (5 points)

Item	Point s	Criteria
Argument 1	2	Explains short initial episodes
Argument 2	2	Explains long mid-training episodes
Argument 3	1	Explains short final episodes

Question 2 (5 points)

Item	Point s	Criteria
Identifies state space	2	States the nature of the state space of Lunar Lander
Explains tabular limitation	2	Explains that tabular Q-learning requires discrete states
Explains DQN advantage	1	Notes that DQN uses a neural network to approximate the Q-function

Optional Qualitative Visualization Bonus

This is optional bonus credit for implementing the qualitative policy visualization

Item	Point s	Criteria
Environment setup for visualization	1	Creates a new LunarLander environment with <code>render_mode="human"</code> or <code>render_mode="rgb_array"</code> and sets agent epsilon to 0 or a low value (≤ 0.1) for evaluation
Evaluation loop implementation	1	Implements a loop that runs multiple episodes (3-10 episodes) with the trained agent, including: environment reset, action selection using the trained agent, environment step, and episode termination handling
Proper evaluation mode	1	Does NOT call <code>agent.update()</code> during visualization (i.e., no learning updates during evaluation) and uses greedy or near-greedy action selection
