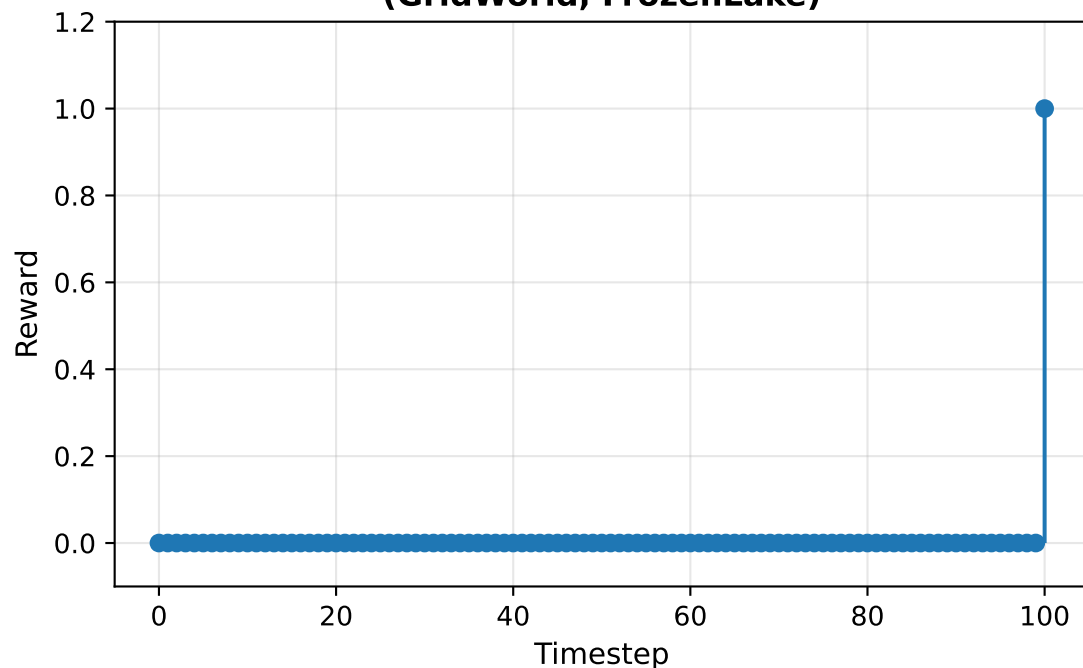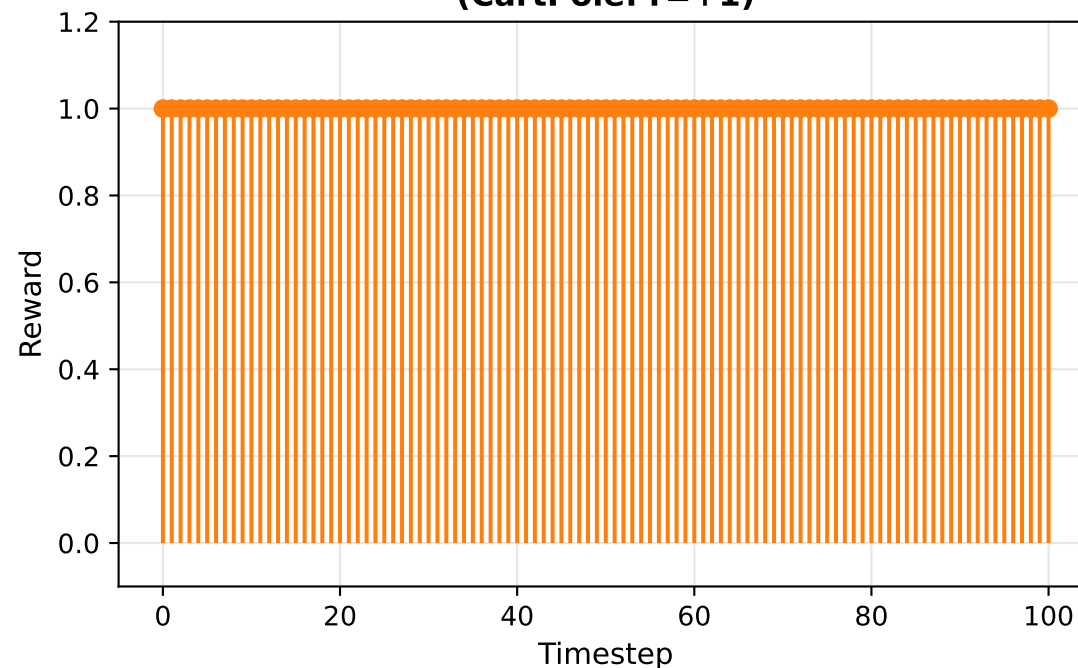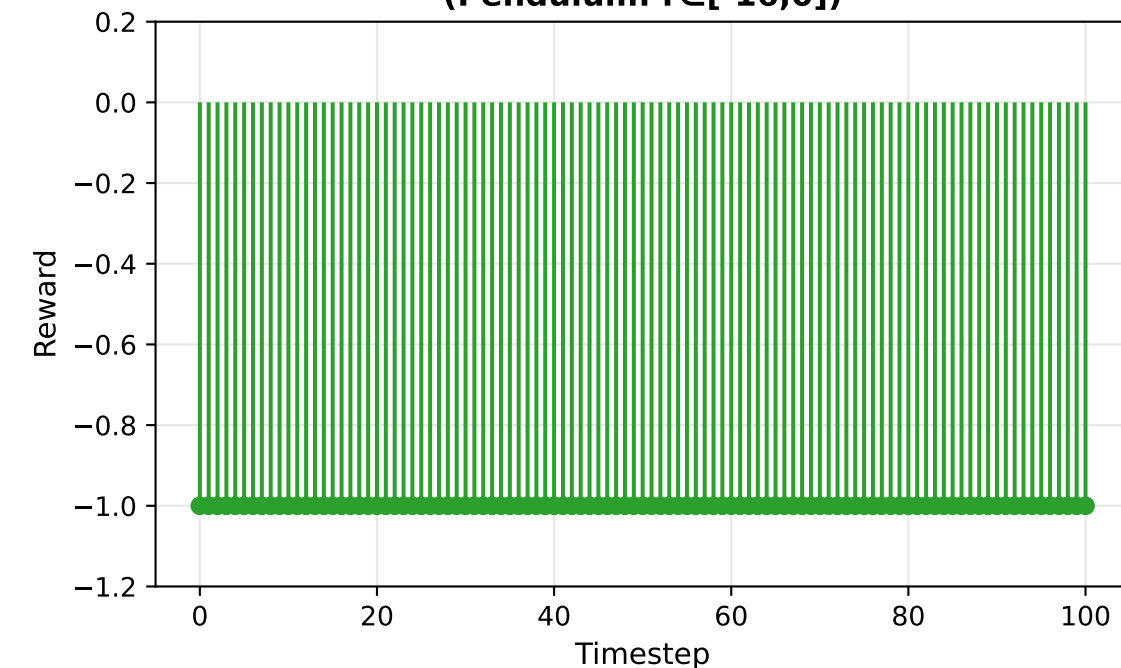**(a) Sparse Reward: Terminal Only (GridWorld, FrozenLake)**
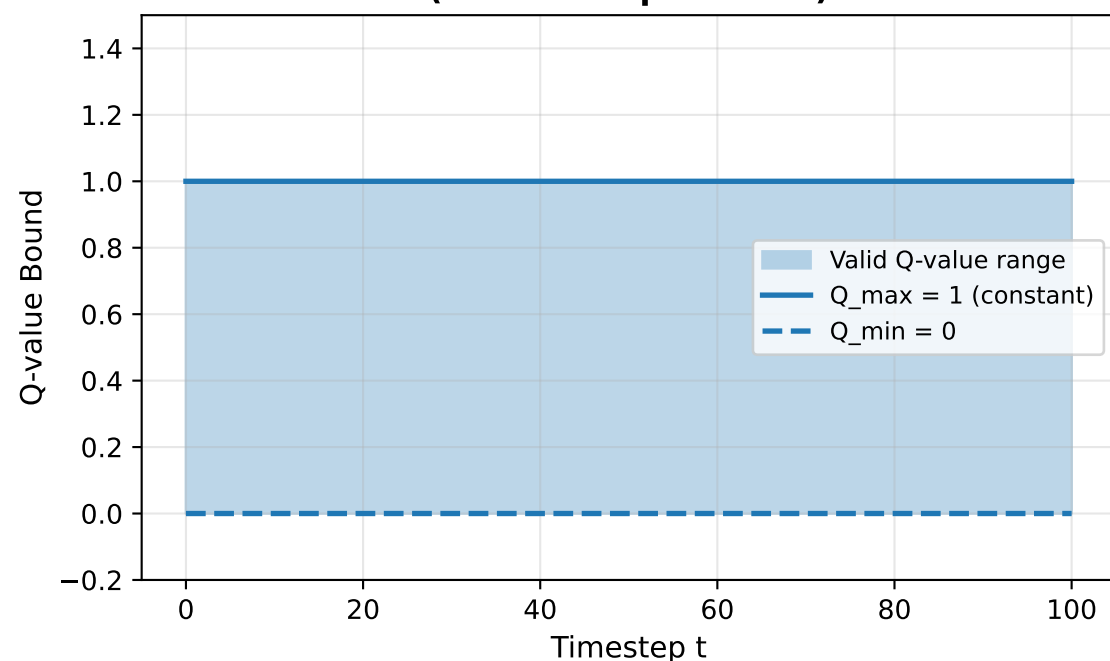
**(b) Dense Positive Reward: Per-Step (CartPole: r=+1)**

**(c) Dense Negative Reward: Per-Step (Pendulum: r∈[-16,0])**
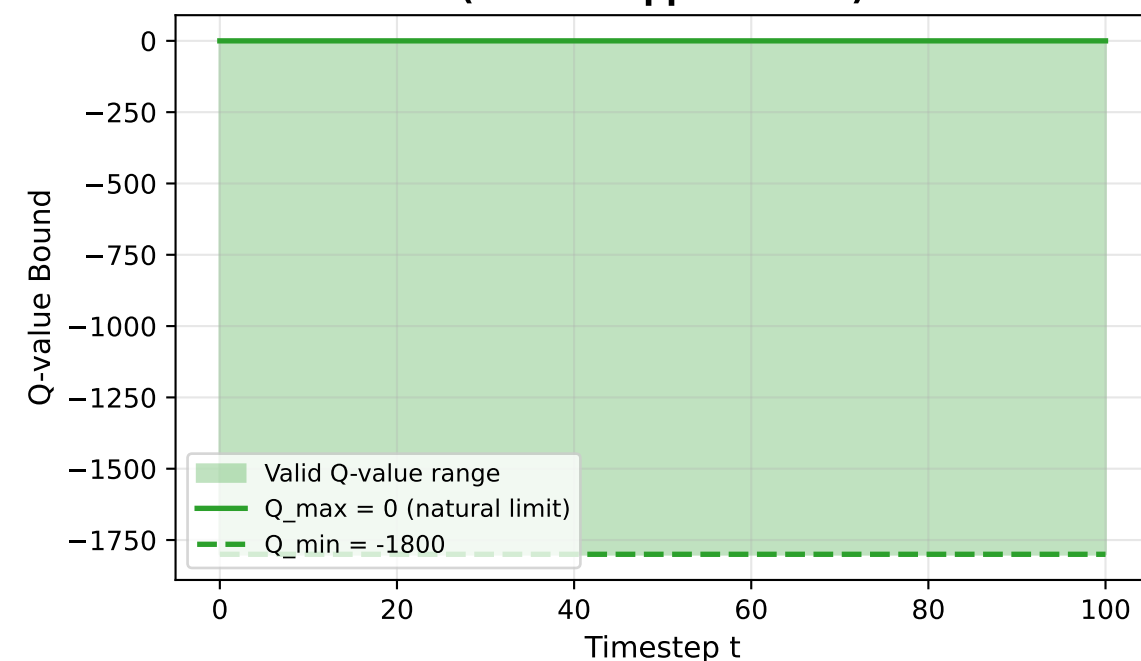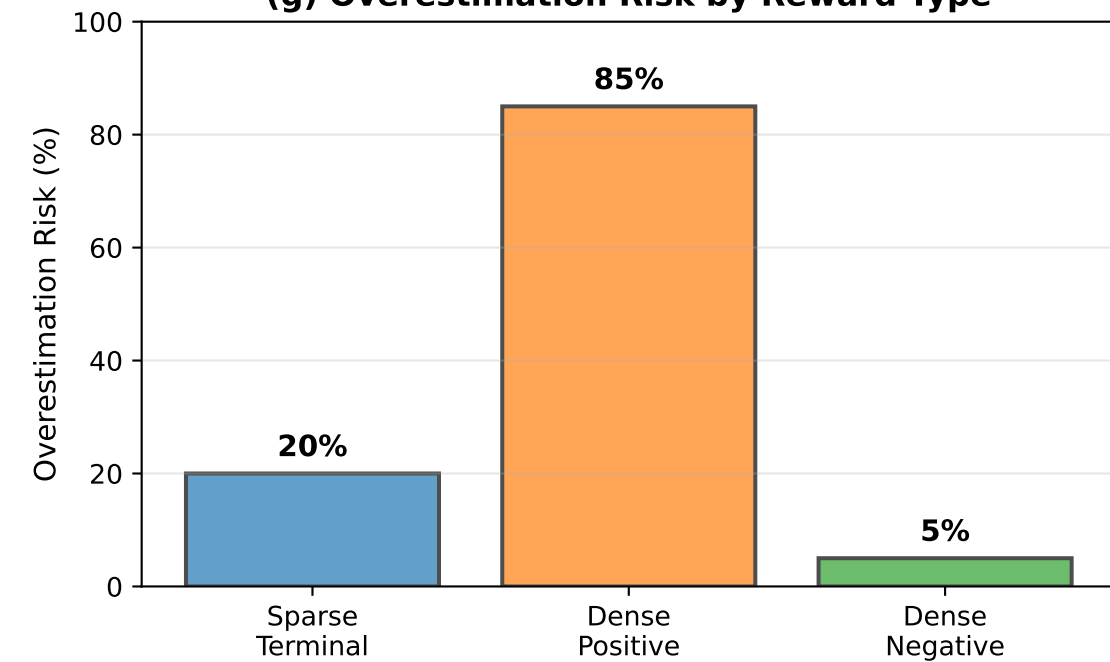
**(d) Sparse: Q-Bounds CONSTANT (No time dependence)**

- Valid Q-value range
- Q_max = 1 (constant)
- Q_min = 0

**(e) Dense Positive: Q_max DECREASES (Remaining potential decreases)**

- Valid Q-value range
- Q_max(t) = (1-γ^(H-t))/(1-γ)
- Q_min = 0

**(f) Dense Negative: Q_max = 0 CONSTANT (Natural upper bound)**

- Valid Q-value range
- Q_max = 0 (natural limit)
- Q_min = -1800

**(g) Overestimation Risk by Reward Type**

- Sparse Terminal: 20%
- Dense Positive: 85%
- Dense Negative: 5%

**(h) Empirical Violation Rates (Without QBound)**

- GridWorld (Sparse): 0.02%
- CartPole (Dense +): 12.5%
- Pendulum (Dense -): 0.0%

**(i) QBound Effectiveness (5 seeds, mean improvement)**

- GridWorld (Sparse): -1.0%
- CartPole (Dense +): +22.5%
- Pendulum (Dense -): -7.0%