

Problems 1 & 2 — Reinforcement Learning Assignment

Problem 1 — Pick-and-Place Robot as an MDP

Task. Learn a policy that controls a robot arm to pick an object from a table and place it at a target quickly and smoothly.

MDP tuple (S, A, P, R, γ) . Episode starts with the object on the table and the arm at a home pose; ends on successful place, failure (collision/drop), or step limit. Discount $\gamma = 0.99$.

1) State S (1 pt). Minimal info for fast, smooth low-level control and task geometry:

- Joint angles/velocities: $q_1..q_n, \dot{q}_1..\dot{q}_n$
- Gripper opening g (0..1)
- End-effector pose: p_{ee} , orientation (quaternion)
- Estimated object pose: p_{obj} , orientation
- Goal (place) pose: p_{goal} , orientation
- Optional phase flag: {approach, grasp, lift, move, place}

2) Action A (1 pt). Continuous control for smoothness:

- Joint torques $\tau_1..\tau_{un}$ (bounded by hardware)
- Gripper command (binary or continuous)

Note: Velocities would also work; torques give more direct smoothness control.

3) Transition P (1 pt). Determined by robot physics + perception updates (contacts, friction, latency). In simulation: physics engine; on the real robot: true plant. Use domain randomization (mass/friction/latency) for robustness.

4) Reward R — speed + smoothness + safety (3 pts). Let $d_{pick}(t)$ = distance $EE \rightarrow$ grasp (pre-grasp); $d_{place}(t)$ = distance object \rightarrow goal (post-grasp); τ_t = joint torques; $\Delta \tau_t = \tau_t - \tau_{t-1}$ (jerk proxy).

- Before grasp: $-\alpha_1 * d_{pick}(t)$
- After grasp: $-\alpha_2 * d_{place}(t)$
- Time cost: $-\lambda$ each step (finish sooner)
- Effort: $-\beta_1 * ||\tau_t||^2$

- Jerk: $-\beta_2 * ||\Delta \tau_t||^2$

Terminal rewards/penalties:

- Success (goal tolerance, gripper open, object stable): $+R_{succ}$

- Collision: $-R_{coll}$

- Dropped object: $-R_{drop}$

- Timeout: $-R_{time}$

Example starting values (tune later): $R_{succ}=100$, $R_{coll}=50$, $R_{drop}=30$, $R_{time}=20$, $\alpha_1=2.0$, $\alpha_2=2.0$, $\lambda=0.1$, $\beta_1=0.001$, $\beta_2=0.0005$. Reasoning: shaping pulls toward grasp/goal; time cost rewards speed; effort/jerk terms enforce low-jerk, smooth, safe motions.

5) Termination (1 pt). Success: object at goal within position/orientation thresholds, gripper released, object static a few steps. Failure: collision, slip/drop, or step limit T .

6) Initial state distribution (1 pt). Randomize object pose in a tray region, small home-pose jitter, and light sensor noise → better generalization.

7) Assumptions & notes (1 pt). Calibrated kinematics; reachable grasp; perception provides pose estimates. Safety layer on robot (action clipping, joint limits, collision monitor). Train with curriculum: approach → grasp → lift → move → place.

Summary. State = robot + task geometry; Actions = continuous torques + gripper; Reward = accuracy + speed + smoothness with clear terminal signals → fast, safe, low-jerk pick-and-place.

Problem 2 — 2×2 Gridworld (Value Iteration, 2 sweeps)

States $S = \{s_1, s_2, s_3, s_4\}$ laid out as: $[s_1 \ s_2; s_3 \ s_4]$. Actions: up, down, left, right. Transitions: deterministic; if a move hits a wall, you stay in place. Rewards: $R(s_1)=5$, $R(s_2)=10$, $R(s_3)=1$, $R(s_4)=2$. Discount $\gamma = 0.9$. Update: $V_{k+1}(s) = R(s) + \gamma \max_a V_k(s')$, where s' is the next state.

Adjacency (next state per action):

State	Next state per action
From s_1	up- $\rightarrow s_1$, left- $\rightarrow s_1$, right- $\rightarrow s_2$, down- $\rightarrow s_3$
From s_2	up- $\rightarrow s_2$, right- $\rightarrow s_2$, left- $\rightarrow s_1$, down- $\rightarrow s_4$
From s_3	down- $\rightarrow s_3$, left- $\rightarrow s_3$, up- $\rightarrow s_1$, right- $\rightarrow s_4$
From s_4	down- $\rightarrow s_4$, right- $\rightarrow s_4$, up- $\rightarrow s_2$, left- $\rightarrow s_3$

Iteration 1

Initial values: $V_0(s_1)=0$, $V_0(s_2)=0$, $V_0(s_3)=0$, $V_0(s_4)=0$

Update to V_1 (all neighbors are 0, so max is 0):

State	V_1 update
-------	--------------

s_1	$5 + 0.9 \cdot 0 = 5$
s_2	$10 + 0.9 \cdot 0 = 10$
s_3	$1 + 0.9 \cdot 0 = 1$
s_4	$2 + 0.9 \cdot 0 = 2$

Greedy actions w.r.t. V_1 (optional): $s_1 \rightarrow \text{right}$; $s_2 \rightarrow \text{up or right}$; $s_3 \rightarrow \text{up}$; $s_4 \rightarrow \text{up}$.

Iteration 2

Use neighbors' V_1 values to compute V_2 :

State	Computation
s1	best among {s1=5, s2=10, s3=1} is 10 => $V2(s1) = 5 + 0.9 \cdot 10 = 14$
s2	best among {s2=10, s1=5, s4=2} is 10 => $V2(s2) = 10 + 0.9 \cdot 10 = 19$
s3	best among {s3=1, s1=5, s4=2} is 5 => $V2(s3) = 1 + 0.9 \cdot 5 = 5.5$
s4	best among {s4=2, s2=10, s3=1} is 10 => $V2(s4) = 2 + 0.9 \cdot 10 = 11$

Final after two iterations: $V2(s1)=14$, $V2(s2)=19$, $V2(s3)=5.5$, $V2(s4)=11$.

Greedy policy w.r.t. $V2$: s1→right; s2→up or right; s3→up; s4→up.