

Reinforcement Learning Notes

Lecture Notes for CSCI 6904 (Deep Reinforcement Learning)

August 19, 2025

Contents

1	Formalization of the RL Problem	6
1.1	Markov Property	6
1.2	Markov Decision Process (MDP)	6
1.3	Goals and Rewards	6
1.3.1	Reward Hypothesis	6
1.3.2	Reward Representation Strategies	6
1.4	Discount Rate γ	7
1.5	Returns and Episodes	7
1.5.1	Problem (Important)	8
1.5.2	Problem (Important)	8
1.6	Partially Observable MDP (POMDP)	9
2	Major Components of an RL Agent	9
2.1	Policy	9
2.2	Value Function	9
2.3	Optimal Value Function	10
2.4	Bellman Expectation Equation	10
2.5	Optimality in the Bellman Equation	10
2.6	Recap	10
2.7	Model	10
2.8	Categorizing RL Agents	11
3	Introduction to Policy Gradient Methods	11
3.1	Policy	11
3.2	Objective	12
4	Policy Gradients for One-Step MDPs / Contextual Bandits	12
4.1	Algorithm (\sim REINFORCE for one-step)	13
4.2	Step-by-Step Derivation of the Policy Gradient (Very important)	13
5	Policy Gradients for Multi-Step MDPs	14
5.1	Objective	14
5.2	Policy Gradient Theorem	14
5.3	REINFORCE Algorithm	14
5.4	REINFORCE with Baseline	14
5.4.1	Algorithm	15

6	Off-Policy Policy Gradients	15
6.1	Concepts	15
6.2	Importance Sampling	15
6.3	Off-Policy REINFORCE with Importance Sampling	15
7	Case Study: Dialog Systems like ChatGPT	16
7.1	Training Phases	16
7.2	Reinforcement Learning with Human Feedback (RLHF)	16
8	Actor-Critic Methods	17
8.1	Recap: Policy Gradient Methods	17
8.1.1	Monte-Carlo Policy Gradient (REINFORCE)	17
8.1.2	REINFORCE with Baseline	17
8.2	Actor-Critic Methods	17
8.2.1	Core Concepts	17
8.2.2	Advantage Actor-Critic	18
8.2.3	(Batch) Advantage Actor-Critic Algorithm	18
8.2.4	Network Designs	18
9	Dynamic Programming (DP)	19
9.1	Bellman Equations	19
9.2	Iterative Policy Evaluation (Prediction)	19
9.3	Policy Iteration (Control)	19
9.4	Value Iteration (Control)	20
10	Asynchronous Dynamic Programming	20
10.1	Synchronous Dynamic Programming	21
10.2	Asynchronous Dynamic Programming	21
10.2.1	In-Place Dynamic Programming	21
10.2.2	Prioritized Sweeping	21
10.2.3	Real-Time Dynamic Programming	21
11	Model-Free Prediction	22
11.1	Monte-Carlo (MC) Learning	22
11.2	Temporal-Difference (TD(0)) Learning	22
11.3	n-Step TD Learning	22
12	Unified View	22
13	TD(λ): Simple Example with $\lambda = 1/2$	23
13.1	Idea	23
13.2	Easy Example	23
14	Model-Free Control	24
14.1	On-Policy Monte Carlo (MC) Control	24
14.1.1	Concept:	24
14.2	Sarsa (On-Policy TD Control)	25
14.3	Q-Learning (Off-Policy TD Control)	26
14.4	n -Step Sarsa	27
14.5	Practical Considerations (All Methods)	27
14.6	Algorithm Boxes (Tabular)	27

14.7 Conceptual Differences Summary	29
15 Deep Q-Networks and DDQN (DQN)	29
15.1 Deep Q-Learning	29
15.2 DQN Algorithm (with Experience Replay and Target Network)	29
15.3 Double DQN: Steps	29
15.4 Prioritized Experience Replay Steps	30
15.5 Improvements to DQN	30
15.5.1 Why N-step Returns Don't Work Simply in DQN	30
16 Continuous Control with Q-Learning	31
16.1 Why DQN Fails for Continuous Actions	31
16.2 Action Discretization	31
16.2.1 Random Shooting Method For DQN	31
16.3 Stochastic Optimization (CEM)	32
16.4 Deep Deterministic Policy Gradient (DDPG)	32
16.4.1 Pseudocode	32
16.5 TD3 - (Not needed for exam)	33
16.5.1 Algorithm (Simplified)	33
16.6 DDPG vs. TD3	33
17 Offline Reinforcement Learning	34
17.1 What is Offline RL?	34
17.2 How is Offline RL Possible?	34
17.3 Why is Offline RL Hard?	34
17.4 Policy Constraint Methods	35
17.5 Conservative Q-Learning (CQL)	35
17.5.1 CQL Algorithm Pseudocode	35
17.6 Additional Methods	35
18 Actor-Critic Methods - Part 3	36
18.1 Batch Advantage Actor-Critic (A2C)	36
18.1.1 Algorithm	36
18.2 PPO - Broken Versions Discussions	36
18.3 Trust Region Policy Optimization (TRPO)	37
18.4 Proximal Policy Optimization (PPO)	38
18.4.1 PPO Algorithm	38
19 Off-Policy Actor-Critic Methods	38
19.1 Online Actor-Critic (On-Policy)	38
19.2 Why Off-Policy Actor-Critic?	39
19.3 Off-Policy Actor-Critic	39
19.4 Deep Deterministic Policy Gradient (DDPG)	39
19.4.1 Pseudocode	40
20 Soft Actor-Critic (SAC)	41
20.1 Maximum Entropy Reinforcement Learning	41
20.2 Soft Value Functions	42
20.3 SAC Algorithm	42

21 Model-Based Reinforcement Learning	42
21.1 Overview	42
21.1.1 Model-Free vs. Model-Based	43
21.2 Advantages and Disadvantages	43
21.2.1 Advantages	43
21.2.2 Disadvantages	43
22 Model Learning	43
22.1 Types of Dynamics Models	43
22.2 Combining Model-Free RL with a Model - Background Planning	43
22.2.1 Model-Based RL: Broken \rightarrow Fixed Progression (MBPO)	44
22.2.2 Dyna-Q Style Algorithm - MBPO	45
22.2.3 Direct Policy Learning with a Model	45
23 Model Based - Decision-Time Planning	46
23.1 Stochastic Optimization	46
23.2 Monte Carlo Tree Search (MCTS)	47
23.3 Model-Based RL with Decision-Time Planning	47
23.4 Uncertainty in RL	47
23.4.1 Sources	47
23.4.2 In Model-Based RL	47
23.4.3 Mitigation	48
23.5 Uncertainty in Neural Networks	48
23.5.1 Estimate model uncertainty (how unsure the parameters are)	48
23.5.2 Bootstrap Ensembles	48
23.5.3 Bayesian Neural Networks (BNNs)	49
23.6 Latent Space Models	49
24 Safe Reinforcement Learning	49
24.1 Overview of Safety Problems	49
24.2 Constrained Markov Decision Processes (CMDP)	50
24.2.1 Machiavelli Benchmark	50
24.2.2 Optimally Safe Policies	50
24.3 Solving CMDPs: Methods	50
24.3.1 Constrained Policy Optimization (CPO)	50
24.3.2 Primal-Dual Methods (Lagrangian)	51
24.3.3 State Augmentation (SauteRL)	52
24.3.4 Action Selection (MASE)	53
24.4 Frontiers in Safe RL	54
25 Exploration in RL	54
25.1 Random Exploration	54
25.2 Novelty Seeking Exploration	54
25.3 Posterior Sampling Methods	54
26 Transfer Learning in RL	55
26.1 Domain Adaptation	55
26.2 Domain Randomization	55
26.3 Multi-Task Transfer	55

27 Frontiers in RL	56
27.1 Meta-Learning	56
27.2 Inverse Reinforcement Learning (IRL)	56
27.3 Hierarchical RL	56
27.4 Foundation Models for RL	56
27.5 Continual RL	56
28 Sample Question and Answer	56

1 Formalization of the RL Problem

1.1 Markov Property

A state S_t is Markov if $p(S_{t+1} | S_t) = p(S_{t+1} | S_1, \dots, S_t)$. This means the future is conditionally independent of the past, given the present state.

Or

The Markov property, also known as the “memoryless property”, is a characteristic of stochastic processes where the future state of a system depends only on its current state, not on its past history.

1.2 Markov Decision Process (MDP)

An MDP is defined as a tuple $\langle \mathcal{S}, \mathcal{A}, P, R, \gamma \rangle$ where:

- \mathcal{S} : The set of states.
- \mathcal{A} : The set of actions.
- P : The state-transition probability matrix, where $p(s'|s, a) = \Pr\{S_t = s' | S_{t-1} = s, A_{t-1} = a\}$.
- R : The reward function, where $r(s, a) = \mathbb{E}[R_t | S_{t-1} = s, A_{t-1} = a]$.
- $\gamma \in [0, 1]$: The discount rate.

Goal: Find the optimal policy $\pi^*(a|s) = \arg \max_{\pi} \mathbb{E}_{\pi}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots | S_t = s]$.

1.3 Goals and Rewards

Rewards define the objective of the problem. A policy’s goal is to maximize the expected cumulative discounted rewards over time.

1.3.1 Reward Hypothesis

“Reward Hypothesis” states that:

“That all of what we mean by goals and purposes can be well thought of as the maximization of the expected value of the cumulative sum of a received scalar signal (called reward).”

This foundational perspective frames reinforcement learning as a way to encode objectives through appropriate design of the reward function.

1.3.2 Reward Representation Strategies

Two common approaches for encoding task objectives in RL are:

- **Goal-Reward Representation:** The agent receives a positive signal (e.g., +1) only upon successfully reaching the goal state; all other transitions yield zero reward.
 - *Pros:* Encourages direct focus on achieving the goal.
 - *Cons:* Rewards are sparse \Rightarrow learning signal is weak; high probability the agent never stumbles upon the goal during exploration.
- **Action-Penalty Representation:** The agent receives a constant negative penalty (e.g., -1) for each non-terminal action, with zero reward at the goal state.

- *Pros*: Provides a denser learning signal pushing the agent toward shorter paths.
- *Cons*: Poorly balanced penalties may discourage exploration; in some cases, the agent may prefer early termination or fail to reach the goal at all if penalties outweigh sparse terminal rewards.

Both schemes highlight the central challenge of reward design: a misspecified objective can lead to ineffective or unintended behavior, despite the reward hypothesis being conceptually sound.

1.4 Discount Rate γ

The discount rate γ controls the preference for immediate versus future rewards. A value of $\gamma = 0$ results in a myopic agent that only considers immediate rewards, while $\gamma = 1$ makes the agent treat all future rewards equally, assuming the task eventually terminates.

1.5 Returns and Episodes

The return G_t is the total discounted reward from time step t . It is defined as:

- For continuing tasks: $G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$.
- For episodic tasks: $G_t = \sum_{k=0}^{T-t-1} \gamma^k R_{t+k+1}$, where T is the terminal time step.

The return can be expressed recursively as $G_t = R_{t+1} + \gamma G_{t+1}$, with $G_T = 0$.

Why the discounted return G_t is finite in continuing tasks

1. What are we worried about? The return is an infinite sum:

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

Since the task does not terminate, we must ensure this sum does not diverge.

2. Why does discounting help? Each later reward is multiplied by γ^k . If $0 \leq \gamma < 1$, the factors $1, \gamma, \gamma^2, \dots$ shrink to 0, so later terms contribute less and less.

3. Bounding the rewards Assume rewards are bounded: $|R_t| \leq R_{\max}$ for all t .

4. Compare with a simpler “worst case” If every reward took its largest magnitude R_{\max} , then

$$G_t \leq R_{\max}(1 + \gamma + \gamma^2 + \gamma^3 + \dots).$$

5. Why does the geometric series converge? For $0 \leq \gamma < 1$,

$$1 + \gamma + \gamma^2 + \dots = \frac{1}{1 - \gamma}.$$

If $\gamma = 1$, the terms never shrink, and the sum goes infinite.

6. Therefore

$$|G_t| \leq \frac{R_{\max}}{1 - \gamma},$$

which is finite. Hence the discounted return G_t is finite in continuing tasks.

Intuition. Discounting makes later rewards smaller; bounded rewards prevent any term from being too large. Summing many shrinking, bounded terms yields a finite result (e.g., $1 + \frac{1}{2} + \frac{1}{4} + \dots = 2$).

1.5.1 Problem (Important)

Question: Given a discount factor $\gamma = 0.8$ and a reward sequence where $R_1 = 5$ and $R_t = 10$ for all $t \geq 2$. Find the return from time 0, G_0 .

Step-by-Step Solution

1. **Define the return formula:** The return from time 0 is given by the discounted sum of all future rewards:

$$G_0 = R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \dots$$

2. **Substitute the given rewards and discount factor:**

$$G_0 = 5 + (0.8)(10) + (0.8)^2(10) + (0.8)^3(10) + \dots$$

3. **Separate the first term and factor out the reward from the remaining terms:**

$$G_0 = 5 + 10 (0.8 + (0.8)^2 + (0.8)^3 + \dots)$$

The terms inside the parentheses form a geometric series.

4. **Sum the geometric series:** The sum of an infinite geometric series of the form $x + x^2 + x^3 + \dots$ where $|x| < 1$ is $\frac{x}{1-x}$. In this case, $x = 0.8$.

$$\text{Sum} = \frac{0.8}{1 - 0.8} = \frac{0.8}{0.2} = 4$$

5. **Calculate the total return:** Substitute the sum of the geometric series back into the equation for G_0 :

$$G_0 = 5 + 10 \times 4 = 5 + 40 = 45$$

1.5.2 Problem (Important)

Given a discount factor $\gamma = 0.8$, a sequence of rewards $R_1 = -3$, $R_2 = 5$, $R_3 = 2$, $R_4 = 7$, and $R_5 = 1$. The terminal time is $T = 5$, which means there are no future rewards after time 5, so $G_5 = 0$. Calculate the return from time 0, G_0 , using the relationship $G_t = R_{t+1} + \gamma G_{t+1}$.

Step-by-Step Solution

1. **Start from the terminal state and work backward:** The terminal return is given as $G_5 = 0$.
2. **Calculate G_4 :** Using the formula $G_4 = R_5 + \gamma G_5$:

$$G_4 = 1 + (0.8)(0) = 1$$

3. **Calculate G_3 :** Using the formula $G_3 = R_4 + \gamma G_4$:

$$G_3 = 7 + (0.8)(1) = 7.8$$

4. **Calculate G_2 :** Using the formula $G_2 = R_3 + \gamma G_3$:

$$G_2 = 2 + (0.8)(7.8) = 2 + 6.24 = 8.24$$

5. **Calculate G_1 :** Using the formula $G_1 = R_2 + \gamma G_2$:

$$G_1 = 5 + (0.8)(8.24) = 5 + 6.592 = 11.592$$

6. **Calculate G_0 :** Using the formula $G_0 = R_1 + \gamma G_1$:

$$G_0 = -3 + (0.8)(11.592) = -3 + 9.2736 = 6.2736$$

1.6 Partially Observable MDP (POMDP)

A POMDP is a generalization of an MDP defined as a tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{O}, P, R, E, \gamma \rangle$, which includes a set of observations \mathcal{O} and an emission probability E , defined as $p(o_t | s_t)$. In a POMDP, the agent's policy is based on observations, $\pi(a|o)$, rather than states.

2 Major Components of an RL Agent

2.1 Policy

An agent's behavior maps states to actions.

- Deterministic: $a = \pi(s)$.
- Stochastic: $\pi(a|s) = p(a|s)$.

2.2 Value Function

Prediction of expected return.

- State-value: $v_\pi(s) = \mathbb{E}_\pi[G_t | S_t = s] = \mathbb{E}_\pi[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s]$.
- Action-value (Q-function): $q_\pi(s, a) = \mathbb{E}_\pi[G_t | S_t = s, A_t = a]$.

If we represent $V_\pi(s)$ in terms of Q value, then:

$$V_\pi(s) = \sum_a \pi(a|s) Q_\pi(s, a)$$

Bellman Equation

The Bellman equation relates a state's value to its immediate reward and the discounted value of the next state.

For a policy π (state-value form):

$$V_\pi(s) = \sum_a \pi(a|s) \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V_\pi(s')]$$

Meaning:

- Weight actions by policy probability.
- Weight next states by transition probability.
- Add immediate reward + discounted future value.

2.3 Optimal Value Function

For optimal $V_*(s)$:

$$V_*(s) = \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V_*(s')]$$

Replace sum over actions with max for best choice.

2.4 Bellman Expectation Equation

The Bellman expectation equation for the state-value function $V_\pi(s)$ is:

$$V_\pi(s) = \mathbb{E}_\pi[R_{t+1} + \gamma V_\pi(S_{t+1}) \mid S_t = s]$$

Here, the expectation operator \mathbb{E}_π represents an average over actions selected by the policy π and the subsequent next states produced by the environment's dynamics. It can be conceptually understood as answering the question: "On average, if I follow policy π from this state, what will I get?"

2.5 Optimality in the Bellman Equation

The relationship between prediction and control problems is central to Reinforcement Learning:

- **Prediction problem:** Given a policy π , find its value function V_π .
- **Control problem:** Find the optimal policy π_* that maximizes the return.

The Bellman optimality equation is what we solve to find the optimal value function V_* :

$$V_*(s) = \max_a \mathbb{E}[R_{t+1} + \gamma V_*(S_{t+1}) \mid S_t = s, A_t = a]$$

Once we have the optimal value function V_* , the optimal policy π_* can be found by acting greedily with respect to it:

$$\pi_*(s) = \arg \max_a \sum_{s'} P(s'|s, a) [R(s, a, s') + \gamma V_*(s')]$$

2.6 Recap

- Bellman expectation: value = average reward + average future value under π .
- Bellman optimality: value = best reward + best future value possible.

2.7 Model

An agent's representation of its environment.

- Dynamics: $p(s', r|s, a) = \Pr\{S_{t+1} = s', R_{t+1} = r \mid S_t = s, A_t = a\}$.
- State transition: $p(s'|s, a)$.
- Reward: $r(s, a) = \mathbb{E}[R_{t+1} \mid S_t = s, A_t = a]$.

2.8 Categorizing RL Agents

- **Policy-based:** A policy, but no value function.
- **Value-based:** A value function, with an implicit policy.
- **Actor-Critic:** Both a policy and a value function.
- **Model-free:** Learns policy or value function, but no model. Learns directly from experience without building model of the environment.
- **Model-based:** A policy or value function, plus a model. Learns the $p(s', r|s, a)$ (state transition probability and the reward) and uses it for planning.

Category	Core Idea	How It Works	Pros	Cons
Policy-based	Learn the policy directly (e.g., REINFORCE , PPO).	Use gradients to adjust policy parameters, often with randomness.	Handles continuous actions well; learns varied strategies.	Can be sample-inefficient and unstable.
Value-based	Learn a value function, pick best actions (e.g., DQN , Double DQN).	Update values (like Q-values) and act greedily.	Efficient for discrete tasks; stable targets.	Struggles with continuous actions; risk of overestimation.
Actor-Critic	Combine policy and value function (e.g., A2C , SAC).	Actor chooses actions, critic evaluates them.	Merges strengths of both; handles continuous actions well.	Coupled training can be unstable; needs careful tuning.
Model-free	Learn from direct experience only (e.g., PPO , DQN).	Learn policy/values via trial and error.	Simple and robust when modeling is hard.	Needs lots of data; no lookahead planning.
Model-based	Learn/use environment model (e.g., Dyna-Q , Dreamer).	Simulate future steps to improve.	Efficient; allows long-term planning.	Model errors can mislead; hard to learn accurate models.

3 Introduction to Policy Gradient Methods

Policy gradient method means it directly optimizes the policy by taking steps in the direction of the gradient of some objective function.

3.1 Policy

The parameterized policy is defined as: $\pi(a|s, \theta) = \Pr(A_t = a \mid S_t = s, \theta_t = \theta)$.

- **Discrete Actions:** A softmax policy, where $\pi(a|s, \theta) = \frac{e^{h(s, a, \theta)}}{\sum_b e^{h(s, b, \theta)}}$, and $h(s, a, \theta)$ is a numerical preference for action a .
- **Continuous Actions:** A Gaussian policy, where the action is sampled from a normal distribution, $a \sim \mathcal{N}(\mu(s, \theta), \sigma^2(s, \theta))$. The probability density function is given by $\pi(a|s, \theta) = \frac{1}{\sigma_\theta(s)\sqrt{2\pi}} \exp\left(-\frac{(a-\mu_\theta(s))^2}{2\sigma_\theta^2(s)}\right)$.

3.2 Objective

The objective is to find the parameter θ that maximizes the expected return $J(\theta)$, which is defined as:

$$J(\theta) \doteq \mathbb{E}_{\pi_\theta}[G_t \mid S_t = s_0] = v_{\pi_\theta}(s_0)$$

Policy gradient methods achieve this by using an iterative update rule, often referred to as **Gradient Ascent**:

$$\theta_{t+1} = \theta_t + \alpha \nabla_\theta J(\theta_t)$$

4 Policy Gradients for One-Step MDPs / Contextual Bandits

One-step MDPs are defined as a process where an agent:

- Starts in a state s sampled from a distribution $d(s)$.
- Takes an action a sampled from the policy $\pi_\theta(a|s)$.
- Receives a reward $r = R_{s,a}$ and the episode terminates.

Example: An online ad system

- States = user contexts: {sports-fan, tech-enthusiast, foodie}.
- The platform sees users with frequencies: $d(\text{sports}) = 0.5$, $d(\text{tech}) = 0.3$, $d(\text{foodie}) = 0.2$ (this is $d(s)$).
- Each time a user arrives:
 - The system observes s sampled from $d(s)$.
 - Chooses an ad (action).
 - Gets a click/no-click reward.
 - The interaction ends.

The objective is to maximize the expected reward, given by:

$$J(\theta) = \mathbb{E}_{\pi_\theta}[r] = \sum_s d(s) \sum_a \pi_\theta(a|s) R_{s,a}$$

The policy gradient is then calculated as:

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(a|s) \cdot r] \approx \frac{1}{N} \sum_i \nabla_\theta \log \pi_\theta(A_i|S_i) \cdot R_{S_i, A_i}$$

The policy parameters are updated using a gradient ascent rule:

$$\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$$

This approach differs from Supervised Learning (SL) Maximum Likelihood in that the RL update is weighted by the reward r , whereas the SL update is an unweighted average.

4.1 Algorithm (\sim REINFORCE for one-step)

1. Run the policy: Sample a state $s \sim d(s)$ and an action $a \sim \pi_\theta(a|s)$.
2. Compute the gradient estimate: $\nabla_\theta J(\theta) = \frac{1}{N} \sum_i \nabla_\theta \log \pi_\theta(A_i|S_i) R_{S_i, A_i}$.
3. Update the policy parameters: $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$.

4.2 Step-by-Step Derivation of the Policy Gradient (Very important)

1. Objective Function: Expected Reward

We want to maximize the expected reward under the current policy:

$$J(\theta) = \mathbb{E}_{\pi_\theta}[r] = \sum_s d(s) \sum_a \pi_\theta(a|s) R_{s,a} \quad (1)$$

Here:

- $d(s)$ is the distribution over states.
- $\pi_\theta(a|s)$ is the probability of taking action a in state s under policy parameters θ .
- $R_{s,a}$ is the reward received when taking action a in state s .

2. Taking the Gradient

We want to compute the gradient $\nabla_\theta J(\theta)$ so we can improve the policy using gradient ascent.

$$\nabla_\theta J(\theta) = \sum_s d(s) \sum_a \nabla_\theta \pi_\theta(a|s) R_{s,a} \quad (2)$$

Note: $d(s)$ and $R_{s,a}$ do not depend on θ , so they are treated as constants.

3. Using the Log Derivative Trick

Start with the chain rule. We want the gradient of a log:

$$\nabla_\theta \log \pi_\theta(a|s) = \frac{1}{\pi_\theta(a|s)} \nabla_\theta \pi_\theta(a|s)$$

Now multiply both sides by $\pi_\theta(a|s)$:

$$\pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) = \nabla_\theta \pi_\theta(a|s)$$

$$\nabla_\theta \pi_\theta(a|s) = \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) \quad (3)$$

Plug this into our gradient:

$$\nabla_\theta J(\theta) = \sum_s d(s) \sum_a \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) R_{s,a} \quad (4)$$

4. Writing as an Expectation

This can now be written as an expectation over samples $(s, a) \sim d(s)\pi_\theta(a|s)$:

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(a|s) \cdot r] \quad (5)$$

Interpretation

This is the policy gradient formula for a one-step MDP:

- It tells us to increase the log-probability of actions that result in high reward.
- This is REINFORCE, simplified to the contextual bandit case.

5 Policy Gradients for Multi-Step MDPs

5.1 Objective

The objective is to maximize the expected return from a starting state, which is equivalent to the state-value function:

$$J(\theta) = \mathbb{E}_{\pi_\theta}[G_t \mid S_t = s_0] = v_{\pi_\theta}(s_0)$$

5.2 Policy Gradient Theorem

The policy gradient theorem provides a way to compute the gradient of the objective function:

$$\nabla_\theta J(\theta) \propto \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(a|s) \cdot q_\pi(s, a)]$$

This is often approximated with the Monte Carlo return G_t :

$$\nabla_\theta J(\theta) \approx \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(a|s) \cdot G_t]$$

This replaces the instantaneous reward with a long-term value $q_\pi(s, a)$ to guide learning.

5.3 REINFORCE Algorithm

1. Sample trajectories $\{\tau_i\}$ from the policy π_θ .
2. For each time step t in a trajectory, compute the return $G_t^i = \sum_{k=t+1}^T \gamma^{k-t-1} R_k^i$.
3. Estimate the policy gradient:

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^{T_i} \nabla_\theta \log \pi_\theta(A_t^i | S_t^i) G_t^i$$

4. Update the policy parameters: $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$.

Note: This algorithm encourages good trajectories by increasing the probability of actions that lead to high returns.

- **Discrete actions:** Use a softmax policy.
- **Continuous actions:** Use a Gaussian policy.

This method suffers from high variance in gradient estimates.

5.4 REINFORCE with Baseline

A baseline $b(s)$ (e.g., the state-value function $v(s)$) can be subtracted from the return to reduce variance without changing the expected value of the gradient.

$$\nabla_\theta J(\theta) \propto \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(a|s) \cdot (q_\pi(s, a) - b(s))]$$

5.4.1 Algorithm

1. Sample trajectories $\{\tau_i\}$ from π_θ .
2. Compute returns G_t^i .
3. Estimate the policy gradient using a learned state-value function $\hat{v}_w(s)$ as the baseline:

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^{T_i} \nabla_\theta \log \pi_\theta(A_t^i | S_t^i) (G_t^i - \hat{v}_w(S_t^i))$$

4. Update the policy parameters: $\theta \leftarrow \theta + \alpha_\theta \nabla_\theta J(\theta)$.
5. Define the value function loss:

$$L(w) = \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^{T_i} (G_t^i - \hat{v}_w(S_t^i))^2$$

6. Update the value function parameters: $w \leftarrow w - \alpha_w \nabla_w L(w)$.

6 Off-Policy Policy Gradients

6.1 Concepts

- **On-policy** methods learn a policy from experience collected by that same policy.
- **Off-policy** methods learn a target policy π using data generated by a different behavior policy μ .

Off-policy learning offers key advantages: it can reuse old data, allows the agent to learn an optimal policy while actively exploring with a different behavior policy, and enables learning about multiple policies simultaneously.

6.2 Importance Sampling

Importance sampling is a technique that allows us to estimate the expected value of a function under one distribution, using samples from another. This is the core of off-policy policy gradients.

$$\mathbb{E}_{x \sim p}[f(x)] = \mathbb{E}_{x \sim q} \left[f(x) \frac{p(x)}{q(x)} \right]$$

The term $\frac{p(x)}{q(x)}$ is the importance sampling weight, which corrects for the difference between the two distributions.

6.3 Off-Policy REINFORCE with Importance Sampling

This method adapts REINFORCE to an off-policy setting by using importance sampling. It allows the algorithm to learn from data collected by an older policy, $\pi_{\theta_{old}}$. The gradient is estimated using the importance sampling weight to correct for the change in policy:

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^{T_i} \frac{\pi_\theta(A_t^i | S_t^i)}{\pi_{\theta_{old}}(A_t^i | S_t^i)} \nabla_\theta \log \pi_\theta(A_t^i | S_t^i) G_t^i$$

The parameters are then updated:

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$$

After the update, the old policy is replaced by the new one:

$$\theta_{old} \leftarrow \theta$$

This approach is notably used in algorithms like Proximal Policy Optimization (PPO), which is essential for training large language models like those in ChatGPT. PPO clips the importance sampling weight to prevent excessively large updates, which helps stabilize learning.

7 Case Study: Dialog Systems like ChatGPT

7.1 Training Phases

The development of advanced dialog systems like ChatGPT typically follows three key phases:

1. **Pretraining (LM)**: A large language model (LM) is first pretrained on a vast corpus of text to learn language patterns and predict the next token.
2. **Instruction Finetuning (SL)**: The pretrained model is then finetuned on a dataset of high-quality instructions and responses using supervised learning (SL) to better follow directions.
3. **Reinforcement Learning with Human Feedback (RLHF)**: The model is further aligned with human preferences using reinforcement learning.

7.2 Reinforcement Learning with Human Feedback (RLHF)

RLHF leverages policy gradients to refine the model's behavior based on human input. The process involves:

1. **Response Sampling**: The policy π_{θ} generates a response (action a) to a given prompt (state s).
2. **Human Feedback**: Human labelers provide feedback on the responses, often by ranking them.
3. **Reward Model Training**: A separate reward model, $r_{\psi}(s, a)$, is trained to predict the human preference score for any response, learning from the collected rankings.
4. **Policy Update**: The policy is updated using a policy gradient algorithm with the reward signal from the reward model. A simplified update rule is:

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log \pi_{\theta}(A_i | S_i) r(S_i, A_i)$$

5. **Over Optimization in Dialog systems**: $J(\theta) = \mathbb{E}_{\pi_{\theta}}[r] - \beta D_{KL}(\pi_{\theta} || \pi_{\beta})$
6. **Off-Policy Considerations**: To improve sample efficiency, off-policy methods with importance sampling are used. A **KL-divergence penalty** is also added to the loss function to prevent the new policy from deviating too drastically from the reference policy, thus maintaining model stability and preventing reward hacking.

8 Actor-Critic Methods

8.1 Recap: Policy Gradient Methods

8.1.1 Monte-Carlo Policy Gradient (REINFORCE)

- **Objective:** The goal is to maximize the expected return from a starting state, represented by the state-value function:

$$J(\theta) = \mathbb{E}_{\pi_\theta}[G_t \mid S_t = s_0] = v_{\pi_\theta}(s_0)$$

- **Policy Gradient Theorem:** The gradient is estimated as the expectation of the log-policy gradient scaled by the action-value function:

$$\nabla_\theta J(\theta) \propto \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(a|s) q_\pi(s, a)]$$

- **Update Rule:** The parameters are updated in the direction of the sampled return G_t , which is a Monte-Carlo estimate of $q_\pi(s, a)$:

$$\theta_{t+1} = \theta_t + \alpha \nabla_\theta \log \pi_\theta(A_t|S_t) G_t$$

This method estimates the action-value function using a complete trajectory's return, which, while unbiased, can have high variance.

8.1.2 REINFORCE with Baseline

To address the high variance of REINFORCE, a baseline $b(s)$ is introduced.

- **Modified Gradient:** The gradient is now proportional to the advantage function, which is the difference between the action-value and a state-dependent baseline. This reduces variance without changing the expected value of the gradient.

$$\nabla_\theta J(\theta) \propto \mathbb{E}_{\pi_\theta}[\nabla_\theta \log \pi_\theta(a|s)(q_\pi(s, a) - b(s))]$$

- **Update Rule:** The update is based on the advantage estimate:

$$\theta_{t+1} = \theta_t + \alpha \nabla_\theta \log \pi_\theta(A_t|S_t)(G_t - b(S_t))$$

- A common and effective baseline is a learned value function, $b(S_t) = \hat{v}(S_t, w)$.

8.2 Actor-Critic Methods

8.2.1 Core Concepts

Actor-Critic methods combine a policy network (the **actor**) with a value function network (the **critic**).

- The **actor** is the policy, parameterized by θ , which chooses the action.
- The **critic** is a value function, parameterized by w , which estimates the value of the state or state-action pair.
- Instead of waiting for a full trajectory's return (Monte Carlo), the actor-critic method uses the critic's value estimate to update the actor. This allows for updates at every step, reducing the variance and speeding up learning.

The critic is used to bootstrap, providing a value estimate to guide the actor.

8.2.2 Advantage Actor-Critic

This approach uses the **advantage function** as the critic's signal to the actor. **The advantage function, $A^\pi(s, a)$, measures how much better a specific action a is compared to the expected value of the state, $V^\pi(s)$.**

- **Advantage function:** $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$

or using the Bellman equation: $A^\pi(s, a) = r + \gamma V^\pi(s') - V^\pi(s)$

- **Advantage estimate:** A common way to estimate the advantage is with the Temporal-Difference (TD) error:

$$\hat{A}(S_t, A_t) = R_{t+1} + \gamma \hat{v}_w(S_{t+1}) - \hat{v}_w(S_t)$$

This is an estimate of how much the reward was better (or worse) than expected.

- **Gradient Update:** The policy gradient is then estimated using this advantage signal:

$$\nabla_\theta J(\theta) \approx \nabla_\theta \log \pi_\theta(A_t|S_t) \hat{A}(S_t, A_t)$$

8.2.3 (Batch) Advantage Actor-Critic Algorithm

This algorithm trains both the actor and the critic networks simultaneously using a batch of sampled trajectories.

1. **Sample Trajectories:** Collect a set of trajectories $\{\tau_i\}$ using the current policy π_θ .
2. **Update Critic:** The critic network's parameters w are updated to minimize the squared TD-error loss.

$$L(w) = \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^{T_i} (R_{t+1}^i + \gamma \hat{v}_w(S_{t+1}^i) - \hat{v}_w(S_t^i))^2$$

The parameters are updated via gradient descent: $w \leftarrow w - \alpha_w \nabla_w L(w)$.

3. **Evaluate Advantage:** For each step in the batch, the advantage estimate is calculated using the updated critic.

$$\hat{A}(S_t^i, A_t^i) = R_{t+1}^i + \gamma \hat{v}_w(S_{t+1}^i) - \hat{v}_w(S_t^i)$$

4. **Update Actor:** The actor network's parameters θ are updated using the batch of advantage estimates.

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^{T_i} \nabla_\theta \log \pi_\theta(A_t^i|S_t^i) \hat{A}(S_t^i, A_t^i)$$

The parameters are updated via gradient ascent: $\theta \leftarrow \theta + \alpha_\theta \nabla_\theta J(\theta)$.

8.2.4 Network Designs

- **Two-Network Design:** The actor and critic are two separate neural networks with their own parameters.
- **Shared-Network Design:** A single network is used, with a shared trunk that processes the state input, and two separate heads for outputting the policy (actor) and the value estimate (critic). This can be more parameter-efficient.

A shared network design is considered parameter efficient because it uses each parameter (a weight in the model) for many different parts of the task.

This approach reduces the total number of parameters required. This, in turn, leads to several benefits for the model, including:

- Faster processing
- Less memory usage
- Better generalization to new tasks

This idea is a core principle in modern deep learning and multi-task learning systems.

9 Dynamic Programming (DP)

9.1 Bellman Equations

- State-value: $v_\pi(s) = \sum_a \pi(a|s) \sum_{s',r} p(s', r|s, a)[r + \gamma v_\pi(s')]$.
- Action-value: $q_\pi(s, a) = \sum_{s',r} p(s', r|s, a)[r + \gamma \sum_{a'} \pi(a'|s') q_\pi(s', a')]$.
- Optimal state-value: $v^*(s) = \max_a \sum_{s',r} p(s', r|s, a)[r + \gamma v^*(s')]$.
- Optimal action-value: $q^*(s, a) = \sum_{s',r} p(s', r|s, a)[r + \gamma \max_{a'} q^*(s', a')]$.

9.2 Iterative Policy Evaluation (Prediction)

“If I follow this strategy, how good is each situation I’ll encounter?” This process does not alter the policy; it simply evaluates it. The algorithm estimates the “value” (long-term expected reward) of each state under a fixed policy π .

Why is this useful? You need to know if a policy is good before improving it. This is the “prediction” part of RL.

Input: An MDP $\langle \mathcal{S}, \mathcal{A}, P, R, \gamma \rangle$ and a policy π .

Output: v_π .

1. Start with an initial guess. Initialize $v_0(s) = 0$ for all $s \in \mathcal{S}$.
2. Repeat until convergence: $v_{k+1}(s) = \sum_a \pi(a|s) \sum_{s',r} p(s', r|s, a)[r + \gamma v_k(s')]$ (Bellman equation).
3. Stop when values stop changing much

9.3 Policy Iteration (Control)

This algorithm finds an optimal policy π^* (the best actions to take in each state) and its corresponding optimal value function V^* .

PolicyIteration = Policy Evaluation + Policy Improvement (This is done iteratively)

It operates by alternating between evaluating the current policy and improving it, a process known as policy iteration.

Why it is useful: Policy iteration is a control method, not just a prediction one. Its goal is to find the best possible strategy to maximize long-term rewards, making it a foundational algorithm for solving control problems in reinforcement learning.

Input: MDP $\langle \mathcal{S}, \mathcal{A}, P, R, \gamma \rangle$.

Output: v^*, π^* .

1. Initialize π_0 arbitrarily (e.g., random actions).
2. Policy Evaluation: Compute v_{π_k} using iterative evaluation (section 9.2).
3. Policy Improvement: $\pi_{k+1}(s) = \arg \max_a q_{\pi_k}(s, a) = \arg \max_a \sum_{s', r} p(s', r | s, a)[r + \gamma v_{\pi_k}(s')]$ (greedy choice based on current v).
4. Repeat until $\pi_{k+1} = \pi_k$.

9.4 Value Iteration (Control)

Similar to Policy Iteration, this algorithm finds an optimal policy π^* and its value function V^* . It combines the evaluation and improvement steps into one by repeatedly applying the Bellman optimality equation to directly compute the optimal value function V^* , from which the optimal policy can be extracted.

Why it is useful Value Iteration is a control method: it not only predicts but also finds the best strategy. It is simpler and often converges faster than Policy Iteration for large problems because it does not require a separate, exhaustive policy evaluation phase – because it combines policy improvement and evaluation into a single update.

1. Initialize $v_0(s) = 0$.
2. Repeat until convergence: $v_{k+1}(s) = \max_a \sum_{s', r} p(s', r | s, a)[r + \gamma v_k(s')]$.
3. Once converged, get Optimal policy: $\pi^*(s) = \arg \max_a \sum_{s', r} p(s', r | s, a)[r + \gamma v^*(s')]$.

10 Asynchronous Dynamic Programming

Asynchronous value iteration is a method for finding the best policy in a Markov Decision Process (MDP) by updating the value of states one at a time, or in small groups, without waiting for a full sweep of all states. This contrasts with traditional, synchronous methods that update all states simultaneously. This flexibility often leads to faster convergence, especially for large state spaces.

How It Works

Imagine a problem of finding the shortest path on a map with many cities (states).

Synchronous Value Iteration would be like a team of people simultaneously calculating the shortest distance from every single city to a destination. They all work on their assigned city at the same time and then, after everyone has finished, they collectively share and update their findings. If the map has a million cities, this is inefficient because the entire team must wait for the slowest person to finish their calculation before the next round of updates can begin.

Asynchronous Value Iteration is like having the same team work independently. One person calculates the distance for City A, another for City B, and so on. As soon as one person finds a new,

better route for their city, they immediately update that city's value. This new information is instantly available for everyone else to use in their ongoing calculations. This continuous and immediate sharing of information allows the process to converge more quickly because no one is idle, waiting for the entire set of calculations to complete.

10.1 Synchronous Dynamic Programming

- All states are **updated in parallel** at each iteration.
- Example: **Value Iteration**.
- For all $s \in S$:

$$v_{k+1}(s) = \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma v_k(s')]$$

- Each $v_{k+1}(s)$ update uses the previous $v_k(s')$ values.

10.2 Asynchronous Dynamic Programming

- States can be **updated in any order**, possibly one at a time.
- Does not require full-sweep of all states per iteration.
- **Converges** as long as all states continue to be selected for updates.

10.2.1 In-Place Dynamic Programming

- Use the latest available values for $v(s)$ during updates, not just those from previous iteration.
- For all $s \in S$:

$$v(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a) [r + \gamma v(s')]$$

10.2.2 Prioritized Sweeping

- Prioritize updates for states with largest Bellman error.
- Leads to faster convergence by focusing computation on relevant parts of the state space.

10.2.3 Real-Time Dynamic Programming

- Updates guided by agent's actual experience.
- Utilizes the environment's feedback to steer which states are updated.
- Helps agent focus learning on states most relevant to its current task.

Pseudocode: Asynchronous Value Iteration

```
for iteration k = 1..K:
    for each selected state s (order may vary):
```

$$v(s) \leftarrow \max_{a,s',r} p(s',r|s,a) [r + v(s')]$$

11 Model-Free Prediction

11.1 Monte-Carlo (MC) Learning

- Learn $v_{\pi}(s) = \mathbb{E}_{\pi}[G_t \mid S_t = s]$.
- MC Update: $V(s) \leftarrow V(s) + \alpha(G_t - V(s))$.
- Incremental: $N(s) \leftarrow N(s) + 1$; Where $N(s)$ is the total number of visits to state s so far (i.e., a state-visitation count)
 $V(s) \leftarrow V(s) + \frac{1}{N(s)}(G_t - V(s))$.

11.2 Temporal-Difference (TD(0)) Learning

- Update: $V(S_t) \leftarrow V(S_t) + \alpha(R_{t+1} + \gamma V(S_{t+1}) - V(S_t))$.
Where :

$$TD.error = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$$

$$Target = (R_{t+1} + \gamma V(S_{t+1}))$$

- TD error: $\delta_t = R_{t+1} + \gamma V(S_{t+1}) - V(S_t)$.

11.3 n-Step TD Learning

- n-step return: $G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n V(S_{t+n})$.
- Update: $V(S_t) \leftarrow V(S_t) + \alpha(G_{t:t+n} - V(S_t))$.

12 Unified View

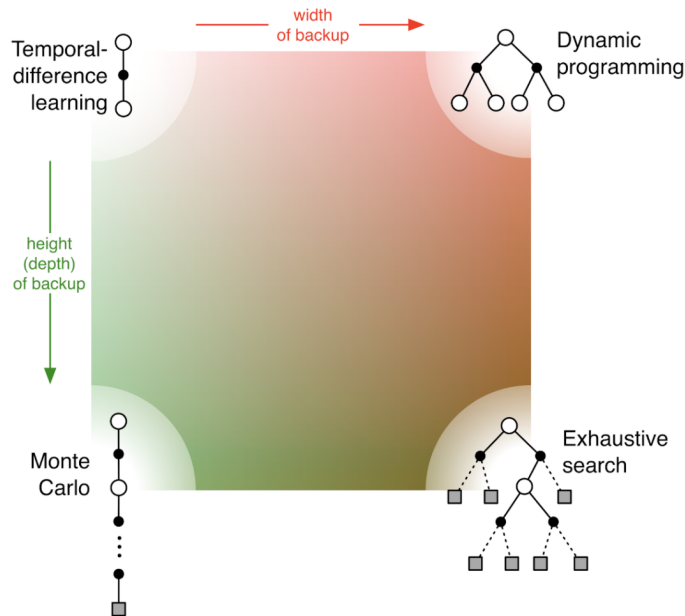


Figure 1: Unified View

What do “depth” and “width” mean?

Depth (length of update): Refers to how far into the future each update looks.

- **Short depth:** Update uses immediate next state (e.g., TD learning).
- **Long depth:** Update could use entire episode or tree (e.g., Monte Carlo, exhaustive search).

Width (branching in update): How many alternative actions or transitions are considered in each update.

- **Narrow width:** Single trajectory or outcome (e.g., real experience).
- **Wide width:** All possible actions/states (e.g., dynamic programming, exhaustive search).

13 TD(λ): Simple Example with $\lambda = 1/2$

13.1 Idea

- TD(λ) averages several n -step returns, with weights λ^{n-1} .
- When $\lambda = 0$: Only uses 1-step TD (TD(0)).
- When $\lambda = 1$: Uses full return (Monte Carlo).
- At $\lambda = \frac{1}{2}$: Each longer return is included but weighted less. The 1-step return gets weight $1/2^0 = 1$, the 2-step return gets weight $1/2^1 = 0.5$, the 3-step return gets weight $1/2^2 = 0.25$, and so on.

$$G_{t:t+n} \doteq R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n \hat{v}(S_{t+n}, \mathbf{w}_{t+n-1}), \quad 0 \leq t \leq T - n,$$

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{T-t-1} \lambda^{n-1} G_{t:t+n} + \lambda^{T-t-1} G_t$$

$$\mathbf{w}_{t+1} \doteq \mathbf{w}_t + \alpha [G_t^\lambda - \hat{v}(S_t, \mathbf{w}_t)] \nabla \hat{v}(S_t, \mathbf{w}_t), \quad t = 0, \dots, T - 1.$$

13.2 Easy Example

Suppose an episode starts in state S_0 , and we observe: $R_1 = 2$, $R_2 = 3$, $R_3 = 4$, $\gamma = 1$ (no discount), and $V(S_3) = 0$ (terminal).

Possible returns:

$$\text{1-step: } G_1^{(1)} = R_1 + V(S_1)$$

$$\text{2-step: } G_1^{(2)} = R_1 + R_2 + V(S_2)$$

$$\text{3-step: } G_1^{(3)} = R_1 + R_2 + R_3 + V(S_3)$$

But let's say $V(S_1) = V(S_2) = V(S_3) = 0$ for simplicity.

So:

$$G_1^{(1)} = 2$$

$$G_1^{(2)} = 2 + 3 = 5$$

$$G_1^{(3)} = 2 + 3 + 4 = 9$$

TD(λ) average for $\lambda = 1/2$:

$$\begin{aligned}
G_1^{(\lambda)} &= (1 - \lambda) \left[G_1^{(1)} + \lambda G_1^{(2)} + \lambda^2 G_1^{(3)} + \dots \right] \\
&= (1 - 0.5) [2 + 0.5 \times 5 + 0.25 \times 9] \\
&= 0.5 [2 + 2.5 + 2.25] \\
&= 0.5 \times 6.75 = 3.375
\end{aligned}$$

Interpretation: With $\lambda = 1/2$, TD(λ) mixes short and long returns: - Gives most weight to 1-step (TD), some to 2-step, less to 3-step, etc. - This offers a balance between quickly updated, low-variance estimates and longer-term, low-bias returns.

Table 1: Summary of TD(λ) returns and weights

<i>n</i> -step Return	Weight (λ^{n-1})	Included Value
1-step	1	2
2-step	0.5	5
3-step	0.25	9

$$\text{TD}(\lambda = 1/2) = 0.5 \times (2 + 2.5 + 2.25) = 3.375$$

This shows how TD(λ) blends multiple returns, with the influence tailing off for longer returns as λ gets smaller.

14 Model-Free Control

We consider a Markov Decision Process (MDP) with a finite state space \mathcal{S} , action space \mathcal{A} , transition kernel $p(s', r \mid s, a)$, and a discount factor $\gamma \in [0, 1)$. The goal of control is to find an optimal policy π^* that maximizes the expected discounted return. All methods below learn an action-value function $Q(s, a)$ directly from experience without an explicit model of the environment's dynamics.

14.1 On-Policy Monte Carlo (MC) Control

Monte Carlo (MC) control improves a policy using returns from complete sampled episodes, following the policy it is evaluating (typically with ϵ -greedy exploration).

$$\pi'(S_t) = \begin{cases} \arg \max_a Q(S_t, a), & \text{with probability } 1 - \epsilon, \\ \text{a random action,} & \text{with probability } \epsilon \end{cases}$$

14.1.1 Concept:

- The agent interacts with the environment using an ϵ -greedy policy for exploration.
- After each episode, it computes the returns G_t for each state-action pair (S_t, A_t) .
- Updates action-value estimates $Q(S_t, A_t)$ using the averages of observed returns.
- The policy is updated to be greedy w.r.t. current Q (or remains ϵ -greedy for continued exploration).

Algorithm Outline:

1. Initialize $Q(s, a)$ arbitrarily, π as ϵ -greedy w.r.t. Q .

2. For each episode:
 - (a) Generate episode following π .
 - (b) For every (S_t, A_t) pair in episode:
 - Compute return G_t from t to episode end.
 - Store all occurrences of (S_t, A_t) .
 - Average all observed returns G_t for (S_t, A_t) :

$$Q(S_t, A_t) \leftarrow \text{average of all returns } G_t \text{ for } (S_t, A_t)$$

- Improve policy: Set π to be ϵ -greedy w.r.t. new Q .

Key Points:

- No bootstrapping—uses actual returns (Monte Carlo).
- Must visit all state-action pairs infinitely often (usually with ϵ -greedy).
- Resulting Q converges to optimal values as episodes accumulate.
- Policy is continually improved as action-value estimates become more accurate.

Control Loop (GLIE in practice)

Repeatedly:

1. Generate an episode by following the current ϵ -greedy policy π .
2. For each visited state-action pair (S_t, A_t) , compute G_t and update Q as described above.
3. Improve π to be ϵ -greedy with respect to the updated Q .

To ensure convergence (in tabular finite MDPs), use GLIE (greedy in the limit with infinite exploration), e.g., $\epsilon_k \downarrow 0$ but $\sum_k \epsilon_k = \infty$. (ϵ will decrease, and at infinity ϵ will be 0, so it will be acting greedy. So to ensure exploration, you need to set epsilon in such a way that sum of all epsilons is equal to ∞).

Notes

- **Pros:** Unbiased targets (true returns), conceptually simple.
- **Cons:** Requires full episodes; high-variance targets; can be slow for long episodes.

14.2 Sarsa (On-Policy TD Control)

Sarsa learns on-policy using one-step bootstrapping.

Update

At each step $(S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1})$:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left(R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t) \right).$$

This is a TD(0) update on the action-value function.

Policy for Exploration

The actions A_t and A_{t+1} are both chosen from the same ϵ -greedy policy π with respect to the Q -function (on-policy).

Control Loop

1. Initialize $Q(s, a)$ arbitrarily; choose ϵ, α .
2. For each episode:
 - (a) Initialize state S_0 , choose action $A_0 \sim \epsilon\text{-greedy}(Q(S_0, \cdot))$.
 - (b) For $t = 0, 1, 2, \dots$ until terminal:
 - Take action A_t , observe reward R_{t+1} and next state S_{t+1} .
 - Choose next action $A_{t+1} \sim \epsilon\text{-greedy}(Q(S_{t+1}, \cdot))$.
 - Update $Q(S_t, A_t)$ with the Sarsa rule.

Notes

- The on-policy nature makes the learned Q reflect exploratory actions, which can be safer (e.g., in a cliff-walking environment).
- Bootstrapped targets reduce variance compared to MC but may introduce bias.

14.3 Q-Learning (Off-Policy TD Control)

Q-learning learns the value of the greedy policy while potentially behaving ϵ -greedily for exploration; it is off-policy via a max target.

Update

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha \left(R_{t+1} + \gamma \max_{a'} Q(S_{t+1}, a') - Q(S_t, A_t) \right).$$

Behavior vs. Target Policies

- **Behavior policy:** The policy used for generating experience (e.g., ϵ -greedy w.r.t. Q).
- **Target policy:** The greedy policy with respect to the current Q -function (implicit in the $\max_{a'}$ operator).

Convergence (Tabular)

Under standard assumptions (finite MDP, sufficient exploration, appropriate stepsizes), Q-learning converges to Q^* with probability 1.

Notes

- It's off-policy and bootstrapped; it can be sensitive to overestimation bias (mitigated by Double Q-learning).
- Often more sample-efficient than MC and is a widely used baseline.

14.4 n -Step Sarsa

n -step Sarsa trades off bias and variance by bootstrapping after n rewards.

Forward n -Step Return

For a transition starting at time t , the truncated return is defined as:

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{n-1} R_{t+n} + \gamma^n Q(S_{t+n}, A_{t+n})$$

with the convention that if the episode terminates before $t + n$, then $Q(S_{t+n}, A_{t+n})$ is omitted (or treated as 0).

Update

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(G_{t:t+n} - Q(S_t, A_t)).$$

Scheduling

At each time t , an update for a prior time step $\tau = t - n + 1$ is performed. This yields a pipeline of delayed updates. As $n \rightarrow 1$, we recover Sarsa; as n approaches the episode length, we approach on-policy MC.

Notes

- Larger n : lower bias, higher variance, and more delay.
- Smaller n : higher bias, lower variance, and faster updates.
- Backward-view equivalents use eligibility traces (e.g., Sarsa(λ)) to implement an exponentially weighted average over different values of n .

14.5 Practical Considerations (All Methods)

- **Initialization:** Optimistic Q_0 values can encourage exploration.
- **Stepsize:** A constant α is good for nonstationary settings, while a diminishing α_t (e.g., Robbins–Monro) supports convergence proofs.
- **Exploration scheduling:** GLIE strategies (e.g., $\epsilon_t \downarrow 0$) balance exploration and exploitation.
- **Terminal handling:** Set Q -values at terminal states to 0.
- **Continuous actions:** The ‘max’ operator in Q-learning is difficult; consider actor-critic or deterministic policy gradient methods instead.

14.6 Algorithm Boxes (Tabular)

On-Policy MC Control (Every-Visit, Incremental)

1. Initialize $Q(s, a)$ arbitrarily; choose $\epsilon \in (0, 1)$, $\alpha \in (0, 1]$.
2. Loop for each episode:
 - (a) Generate episode $(S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T)$ using an ϵ -greedy policy derived from Q .

(b) For $t = 0, \dots, T - 1$:

$$G_t \leftarrow \sum_{k=0}^{T-t-1} \gamma^k R_{t+1+k}$$

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(G_t - Q(S_t, A_t)).$$

Sarsa (On-Policy TD(0) Control)

1. Initialize $Q(s, a)$ arbitrarily; choose ϵ, α .
2. Loop for each episode:
 - (a) Initialize state S_0 , choose action $A_0 \sim \epsilon\text{-greedy}(Q(S_0, \cdot))$.
 - (b) For $t = 0, 1, 2, \dots$ until terminal:
 - Take A_t , observe R_{t+1}, S_{t+1} .
 - Choose $A_{t+1} \sim \epsilon\text{-greedy}(Q(S_{t+1}, \cdot))$.
 - Update $Q(S_t, A_t)$ using the Sarsa rule:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)).$$

Q-Learning (Off-Policy TD Control)

1. Initialize $Q(s, a)$ arbitrarily; choose ϵ, α .
2. Loop for each episode:
 - (a) Initialize state S_0 .
 - (b) For $t = 0, 1, 2, \dots$ until terminal:
 - Choose $A_t \sim \epsilon\text{-greedy}(Q(S_t, \cdot))$.
 - Take A_t , observe R_{t+1}, S_{t+1} .
 - Update $Q(S_t, A_t)$ using the Q-Learning rule:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha(R_{t+1} + \gamma \max_{a'} Q(S_{t+1}, a') - Q(S_t, A_t)).$$

n -Step Sarsa (Forward View)

1. Initialize $Q(s, a)$ arbitrarily; choose $\epsilon, \alpha, n \geq 1$.
2. Loop for each episode:
 - (a) Generate a trajectory $(S_0, A_0, R_1, S_1, A_1, \dots)$ using an ϵ -greedy policy derived from Q .
 - (b) For $t = 0, 1, 2, \dots$ until all updates applied:
 - Let $\tau \leftarrow t - n + 1$.
 - If $\tau \geq 0$, define the n -step return:

$$G_{\tau:\tau+n} = \sum_{k=0}^{n-1} \gamma^k R_{\tau+1+k} + \gamma^n Q(S_{\tau+n}, A_{\tau+n})$$

(with truncation at terminal states), and update:

$$Q(S_\tau, A_\tau) \leftarrow Q(S_\tau, A_\tau) + \alpha(G_{\tau:\tau+n} - Q(S_\tau, A_\tau)).$$

14.7 Conceptual Differences Summary

- **Target:** MC uses the full return G_t ; Sarsa uses a one-step bootstrapped target $R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$; Q-learning uses a max-operator bootstrapped target $R_{t+1} + \gamma \max_{a'} Q(S_{t+1}, a')$; n -step Sarsa uses an intermediate n -step return $G_{t:t+n}$.
- **On- vs Off-Policy:** MC control and Sarsa are on-policy; Q-learning is off-policy.
- **Bias/Variance:** MC has low bias but high variance; TD methods have more bias but lower variance. The n -step methods provide a spectrum to balance this trade-off by adjusting the step size n .

15 Deep Q-Networks and DDQN (DQN)

15.1 Deep Q-Learning

- Represent $Q(s, a; w) \approx q^*(s, a)$.
- Loss: $L(w) = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i; w))^2$, where $y_i = r_i + \gamma \max_{a'} Q(s'_i, a'; w^-)$.
- Update w via stochastic gradient descent (SGD).

15.2 DQN Algorithm (with Experience Replay and Target Network)

1. Initialize replay buffer \mathcal{B} ; initialize online Q-network $Q(s, a; w)$; initialize target network $Q(s, a; w^-) \leftarrow Q(s, a; w)$.
2. For episode = 1, ...:
 - (a) Reset env, receive initial state s_0 ; set $t \leftarrow 0$.
 - (b) While not done:
 - i. Select a_t via ϵ -greedy from $Q(\cdot; w)$ given s_t .
 - ii. Execute a_t , observe r_{t+1} , next state s_{t+1} , and done flag $d_{t+1} \in \{0, 1\}$.
 - iii. Store $(s_t, a_t, r_{t+1}, s_{t+1}, d_{t+1})$ in \mathcal{B} .
 - iv. Sample a mini-batch $\{(s_i, a_i, r_i, s'_i, d_i)\}_{i=1}^N$ from \mathcal{B} .
 - v. Compute targets: $y_i = r_i + \gamma(1 - d_i) \max_{a'} Q(s'_i, a'; w^-)$.
 - vi. Update w by minimizing $L(w) = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i; w))^2$ (often Huber loss).
 - vii. Every N steps: update target, e.g., hard $w^- \leftarrow w$ or Every K steps soft $w^- \leftarrow \tau w^- + (1 - \tau)w$.
 - viii. Optionally anneal ϵ ; $t \leftarrow t + 1$; $s_t \leftarrow s_{t+1}$.

15.3 Double DQN: Steps

1. **Store transition:** Take action A_t , store $(S_t, A_t, R_{t+1}, S_{t+1})$ in the replay buffer B .
2. **Sample mini-batch:** Randomly pick a mini-batch of transitions $\{(S, A, R, S')\}$ from B .
3. **Calculate target:**
 - **DQN target:** $y_i = R_{i+1} + \gamma \hat{q}_{w^-}(S_{i+1}, \arg \max_{a'} \hat{q}_w(S_{i+1}, a'))$
 - **Double DQN target:** $y_i = R_{i+1} + \gamma \hat{q}_{w^-}(S_{i+1}, \arg \max_{a'} \hat{q}_w(S_{i+1}, a'))$

4. **Calculate loss:** $L(w) = \sum_i (y_i - \hat{q}(S_i, A_i, w))^2$
5. **Update weights:** Adjust w to minimize $L(w)$.
6. **Target network update:** Set $w^- \leftarrow w$ every N steps.

Alternative target network update: $w^- \leftarrow \tau w^- + (1 - \tau)w$, e.g., $\tau = 0.999$ every step.

15.4 Prioritized Experience Replay Steps

1. **Store transition:** Take action A_t , and store $(S_t, A_t, R_{t+1}, S_{t+1})$ in replay buffer B with maximal priority $p_t = \max_i p_i$.
2. **Sample mini-batch:** Sample transitions from B using prioritized experience replay: probability $P(j) = \frac{p_j^\alpha}{\sum_i p_i^\alpha}$
3. **Importance-sampling:** Compute weight $\lambda_i = \left(\frac{1}{NP(i)} \right)^\beta / \max_i \lambda_i$
4. **Update priority:** Set $p_j \leftarrow |\delta_j| + \epsilon$ for sampled transitions.
5. **Calculate target:** $y_i = R_{i+1} + \gamma \hat{w}^-(S_{i+1}, \arg \max_{a'} \hat{w}(S_{i+1}, a'))$
6. **Calculate loss:** $L(w) = \sum_i \lambda_i (y_i - \hat{q}(S_i, A_i, w))^2$
7. **Update w :** Minimize $L(w)$ and update network parameters.
8. **Update target network:** $w^- \leftarrow \tau w^- + (1 - \tau)w$ every step.

15.5 Improvements to DQN

- **Double DQN:** Target $y_i = r_i + \gamma Q(s'_i, \arg \max_{a'} Q(s'_i, a'; w), w^-)$ to reduce overestimation.
- **Prioritized Replay:** Sample transitions with priority $p_i \propto |\delta_i| + \epsilon$; update with importance sampling weights.
- **Multi-step Returns:** $y_i = \sum_{k=1}^n \gamma^{k-1} r_{i+k} + \gamma^n \max_{a'} Q(s'_{i+n}, a'; w^-)$.

15.5.1 Why N-step Returns Don't Work Simply in DQN

Only actually correct when learning on-policy. Let's assume the experience replay buffer (ERB) contains a sequence of transitions:

$$\{(S_t, A_t, R_{t+1}, S_{t+1}, A_{t+1}, R_{t+2}, S_{t+2}), \dots\}$$

This sequence was collected using an ϵ -greedy behavior policy.

For a single-step return, the tuple $(S_t, A_t, R_{t+1}, S_{t+1})$ can be safely used. The reward R_{t+1} is determined entirely by the environment and does not depend on the specific Q-network that generated the action A_t .

However, for an N-step return, we also need the subsequent actions, such as A_{t+1} , from state S_{t+1} . In standard Deep Q-Networks (DQN), the target is computed using the target network, specifically with $\max_{a'} Q_{\text{target}}(S', a')$, not the action actually taken by the behavior policy.

If we directly use A_{t+1} from the Experience Replay Buffer, we are implicitly assuming that it was chosen according to the target network's policy. In reality, A_{t+1} was selected by the older behavior policy.

Since the Q-network parameters may have changed significantly since that transition was collected, **this assumption is incorrect and introduces an off-policy bias into the N-step return calculation.**

Simple example:

- Step 1: Agent uses policy π_1 (old) to generate transition
- Step 2: Agent uses policy π_2 (newer) to generate next transition
- Step 3: We sample both from replay buffer and compute n-step return
- **Problem:** N-step assumes single policy, but we mixed π_1 and π_2 !

Solution:

- Ignore the problem – often works very well for DQN.
- Collect transitions $\{S, A, R, S', A', R', S'', \dots\}$ and keep the n-step trace only if the target network's action for S' (and for any subsequent states) matches the recorded actions; otherwise, cut the trace.
- Importance sampling

Key insight: Actor-critic learns from fresh on-policy data. DQN learns from mixed old/new off-policy data stored in replay buffer.

16 Continuous Control with Q-Learning

16.1 Why DQN Fails for Continuous Actions

Problem: DQN is designed for **discrete action spaces**.

1. **Q-value maximization:** DQN requires computing $\max_a Q(s, a)$ over all possible actions a . For continuous actions, this would mean searching an *infinite* number of actions—impractical to compute or represent exactly.
2. **Action selection:** The neural network output must be one value per action. In continuous control, there are infinitely many actions, so mapping actions to outputs is not feasible.
3. **Argmax not tractable:** For continuous actions, finding the action a that maximizes $Q(s, a)$ is a complex optimization problem, not just selecting the highest output.

Conclusion: DQN cannot be directly applied for environments where the action space is continuous. Instead, algorithms like DDPG or SAC are used for such cases.

16.2 Action Discretization

- Discretize continuous actions into bins (e.g., 10 bins per dimension).

16.2.1 Random Shooting Method For DQN

- Randomly sample N possible continuous actions $\{a_1, a_2, \dots, a_N\}$.
- Evaluate each action using Q-function: $\hat{q}_w(s, a_i)$.
- Pick the action with the highest value: $\max_i \hat{q}_w(s, a_i) \varepsilon$ greedy times and put them into the Experience Replay Buffer.
- As N increases, the approximation to the optimal action improves.

16.3 Stochastic Optimization (CEM)

1. Sample M actions from distribution $p(a)$ (e.g., Gaussian $\mathcal{N}(\mu, \sigma^2)$).
2. Evaluate $Q(s, a_i)$ for each.
3. Select top K elites.
4. Refit $p(a)$ to elites (e.g., update μ, σ).
5. Best action: $\arg \max_a Q(s, a) \approx$ elite mean or max.

16.4 Deep Deterministic Policy Gradient (DDPG)

- DDPG is an **actor-critic** method for handling **continuous action spaces**.
- The **actor network** μ_θ outputs the best action for a given state S . It is called the critic network because now the objective is to maximize the critic.
- The **critic network** $Q_w(S, A)$ estimates the action-value function.
- The network has two parts: an **actor** $\mu(s)$ and a **critic** $Q(s, a)$.
- The “deterministic” in DDPG refers to the type of policy the algorithm learns. A deterministic policy outputs a specific action for each given state, as opposed to a stochastic policy which outputs a probability distribution over actions.

16.4.1 Pseudocode

1. Initialize actor and critic networks μ and Q with parameters θ and w .
2. Initialize target networks μ^- and Q^- with parameters $\theta^- \leftarrow \theta$ and $w^- \leftarrow w$.
3. Initialize replay buffer \mathcal{B} .
4. For each episode:
 - (a) Initialize a random process \mathcal{N} for exploration.
 - (b) Receive initial state S .
 - (c) For $t = 1, T$:
 - i. Select action $A = \mu_\theta(S) + \mathcal{N}_t$ (Noise).
 - ii. Execute action A and observe new state S' and reward R .
 - iii. Store transition (S, A, R, S') in \mathcal{B} .
 - iv. Sample a random minibatch of N transitions from \mathcal{B} .
 - v. For each sampled transition (S_i, A_i, R_i, S'_i) , calculate the target value:

$$y_i = R_i + \gamma Q_{w^-}(S'_i, \mu_{\theta^-}(S'_i))$$

- vi. Update the critic network by minimizing the loss:

$$L(w) = \frac{1}{N} \sum_i (y_i - Q_w(S_i, A_i))^2$$

vii. Update the actor network using the sampled policy gradient:

$$\nabla_{\theta} J \approx \frac{1}{N} \sum_i \nabla_A Q_w(S_i, A)|_{A=\mu_{\theta}(S_i)} \nabla_{\theta} \mu_{\theta}(S_i)$$

– Now the objective is to maximize the critic.

viii. Soft update the target networks:

$$w^- \leftarrow \tau w + (1 - \tau) w^-$$

$$\theta^- \leftarrow \tau \theta + (1 - \tau) \theta^-$$

16.5 TD3 - (Not needed for exam)

Twin Delayed Deep Deterministic Policy Gradient

Key Idea: TD3 = DDPG + 3 stability fixes to reduce overestimation bias.

1. **Twin Critics:** Use $Q_{\theta_1}, Q_{\theta_2}$ and take:

$$y = r + \gamma \min_{i=1,2} Q_{\theta'_i}(s', a')$$

to avoid optimistic bias.

2. **Target Policy Smoothing:** Add small clipped noise to target action:

$$a' = \pi_{\phi'}(s') + \epsilon, \quad \epsilon \sim \text{clip}(\mathcal{N}(0, \sigma), -c, c)$$

for smoother policies.

3. **Delayed Policy Updates:** Update actor & targets every d steps, letting critics learn first.

16.5.1 Algorithm (Simplified)

1. Init $Q_{\theta_1}, Q_{\theta_2}, \pi_{\phi}$, target networks $Q_{\theta'_1}, Q_{\theta'_2}, \pi_{\phi'}$, and replay buffer \mathcal{B} .

2. For each step:

(a) Execute $a_t = \pi_{\phi}(s_t) + \text{exploration noise}$, store (s_t, a_t, r_t, s_{t+1}) in \mathcal{B} .

(b) Sample minibatch from \mathcal{B} .

(c) Compute target:

$$y = r + \gamma \min_{i=1,2} Q_{\theta'_i}(s', \pi_{\phi'}(s') + \epsilon)$$

(d) Update critics by MSE to y .

(e) If $t \bmod d = 0$:

$$\begin{aligned} \nabla_{\phi} J(\phi) &= \frac{1}{N} \sum Q_{\theta_1}(s, \pi_{\phi}(s)) \\ \theta'_i &\leftarrow \tau \theta_i + (1 - \tau) \theta'_i, \quad \phi' \leftarrow \tau \phi + (1 - \tau) \phi' \end{aligned}$$

16.6 DDPG vs. TD3

Feature	DDPG	TD3
Critics	1	2 (take min)
Target smoothing	No	Yes
Actor update frequency	Every step	Delayed (e.g. every 2)

17 Offline Reinforcement Learning

17.1 What is Offline RL?

Offline RL (also known as batch RL or fully off-policy RL) is a subfield of reinforcement learning where the goal is to learn an optimal policy π from a fixed, pre-collected dataset $D = \{(s_i, a_i, s'_i, r_i)\}$ without any further interaction with the environment. The dataset can be collected by any behavior policy π_β (e.g., a random policy, an expert human, or a series of past RL runs).

The objective is to find a policy π that maximizes the expected cumulative reward:

$$\max_{\pi} \sum_{t=0}^T \mathbb{E}_{s_t \sim d^{\pi}, a_t \sim \pi(a|s)} [r(s_t, a_t)]$$

Once learned, the policy is deployed in the real environment. A related task is **Off-Policy Evaluation (OPE)**, which aims to estimate the value of a given policy π from the offline dataset D , i.e., $J(\pi) = \mathbb{E}_{\pi}[\sum_{t=1}^T r(s_t, a_t)]$.

17.2 How is Offline RL Possible?

Offline RL can be effective by:

- **Extracting good behaviors:** Learning to distinguish between high-reward and low-reward behaviors within a mixed-quality dataset.
- **Generalization:** A good action observed in one state can inform the model about good actions in similar, unseen states.
- **Stitching:** The algorithm can recombine segments of different trajectories from the dataset to form a new, superior policy. This is sometimes more powerful than imitation learning, as it can create behaviors that were never seen in their entirety.

A key example is offline QT-Opt for robotic grasping, which achieved a high success rate by learning from a large, fixed dataset of grasping attempts.

17.3 Why is Offline RL Hard?

The primary challenge in offline RL is the **distribution shift** between the learned policy π and the data-collecting behavior policy π_β .

- **Counterfactual queries:** The learned policy may select **out-of-distribution (OOD)** actions that are not present in the dataset. Since the value function is only trained on in-distribution actions, it may produce wildly inaccurate and overly optimistic Q-value estimates for these OOD actions.
- **Generalization issues:** Function approximation, especially with deep neural networks, can amplify these errors, leading to a policy that learns to exploit these erroneous Q-values and performs poorly in the real environment. This is a common failure mode, where naive applications of off-policy algorithms like DQN can lead to a significant overestimation of Q-values on offline data.

17.4 Policy Constraint Methods

One way to mitigate distribution shift is to constrain the learned policy to stay close to the behavior policy. This can be formalized as:

$$\pi(a|s) = \arg \max_{\pi} \mathbb{E}[Q(s, a)] \text{ s.t. } D_{KL}(\pi || \pi_{\beta}) \leq \epsilon$$

This approach can be difficult because the behavior policy π_{β} is often unknown. Moreover, if the constraint is too strict, the policy may not be able to improve, and if it's too loose, the distribution shift problem persists.

17.5 Conservative Q-Learning (CQL)

Conservative Q-Learning (CQL) is a method that directly addresses the overestimation problem by modifying the Q-learning objective. Its goal is to **conservatively lower the Q-values for OOD actions** while keeping the Q-values for in-distribution actions high. The CQL loss function is:

$$L_{CQL}(Q) = \alpha \left(\mathbb{E}_{s \sim D, a \sim \mu(a|s)}[Q(s, a)] - \mathbb{E}_{(s, a) \sim D}[Q(s, a)] \right) + \mathbb{E}_{(s, a, s') \sim D} \left[\left(Q(s, a) - (r(s, a) + \gamma \mathbb{E}_{a' \sim \pi(a'|s')}[Q(s', a')]) \right)^2 \right]$$

- Where μ is the sampling distribution over actions (Softmax).

The first term is the key addition: it pushes down the Q-values on a set of sampled OOD actions (e.g., from a uniform distribution μ) while pushing up the Q-values for actions present in the dataset.

17.5.1 CQL Algorithm Pseudocode

1. Initialize the Q-network $Q(s, a; w)$, policy network $\pi(a|s; \theta)$, and populate a replay buffer with the offline dataset D .
2. For each training step:
3. Sample a mini-batch of transitions from D .
4. Update the Q-network parameters w using a stochastic gradient descent (SGD) step on the $L_{CQL}(Q)$ loss.
5. Update the policy parameters θ :
 - **For discrete actions:** The policy is updated to be greedy with respect to the learned Q-values: $\pi(a|s) = \arg \max_{a'} Q(s, a')$.
 - **For continuous actions:** The policy is updated to maximize the expected Q-value: $\theta \leftarrow \theta + \alpha \nabla_{\theta} \mathbb{E}_{a \sim \pi(a|s)}[Q(s, a)]$.

17.6 Additional Methods

- **Implicit Q-Learning (IQL):** A method that avoids explicitly querying OOD actions by using an expectile regression loss.
- **Advantage-Weighted Actor-Critic (AWAC):** An algorithm that performs policy updates by weighting offline samples based on their estimated advantage.

- **Model-Based Methods:** Algorithms like MOPO (Model-based Offline Policy Optimization) and COMBO (Conservative Model-Based Optimization) learn a world model from the offline data and use it to generate synthetic transitions, while often incorporating a conservative penalty to avoid model-exploitation.
- **Sequence Modeling:** Methods like the Trajectory Transformer treat RL as a sequence modeling problem, learning to predict future returns and actions from past trajectories.

18 Actor-Critic Methods - Part 3

18.1 Batch Advantage Actor-Critic (A2C)

A2C is a stable on-policy method that reduces variance compared to REINFORCE by using a learned value function to estimate the advantage.

- **Advantage function:** The advantage is estimated using the TD-error:

$$\hat{A}(S_t, A_t) = R_{t+1} + \gamma \hat{v}_w(S_{t+1}) - \hat{v}_w(S_t)$$

18.1.1 Algorithm

1. Sample a batch of trajectories $\{\tau_i\}$ using the current policy $\pi_\theta(a|s)$.
2. **Critic Update:** Update the critic network parameters w to minimize the Mean Squared Error (MSE) loss on the TD-error:

$$L(w) = \frac{1}{N} \sum_i \sum_t (R_{t+1}^i + \gamma \hat{v}_w(S_{t+1}^i) - \hat{v}_w(S_t^i))^2$$

$$w \leftarrow w - \alpha_w \nabla_w L(w)$$

3. **Advantage Evaluation:** Calculate the advantage estimate for each time step in the batch.
4. **Actor Update:** Update the actor network parameters θ using the advantage-weighted policy gradient:

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i \sum_t \nabla_\theta \log \pi_\theta(A_t^i | S_t^i) \hat{A}(S_t^i, A_t^i)$$

$$\theta \leftarrow \theta + \alpha_\theta \nabla_\theta J(\theta)$$

18.2 PPO - Broken Versions Discussions

1. PPO v0.5 — Broken

What it does:

- Runs like a standard Advantage Actor-Critic (A2C) update but just uses new parameters θ every iteration:

$$\nabla J(\theta) \approx \frac{1}{N} \sum \nabla_\theta \log \pi_\theta(a_i | s_i) \hat{A}(s_i, a_i)$$

- Data is collected with $\pi_{\theta_{\text{old}}}$, but the log-probability ratio $\frac{\pi_\theta}{\pi_{\theta_{\text{old}}}}$ is not used — it ignores that the policy changed between data collection and gradient computation.

Why is this broken:

- Large gradient steps can move π_θ far from $\pi_{\theta_{\text{old}}}$ in one update.
- On-policy assumption is violated because the collected batch no longer reflects the updated policy.
- Can cause destructive policy updates (performance collapse after one bad step).

2. PPO v0.8 — Still Broken (But Closer)

What it adds:

- Introduces the ratio term:

$$J_{\text{PPO}}(\theta) = \mathbb{E}_{\pi_{\theta_{\text{old}}}} \left[\frac{\pi_\theta(a|s)}{\pi_{\theta_{\text{old}}}(a|s)} \hat{A}_{\pi_{\theta_{\text{old}}}}(s, a) \right]$$

- This is the importance-sampled surrogate objective from TRPO theory.

Why still broken:

- The ratio is used, but there's no constraint or clipping — the ratio can still blow up.
- High variance and instability remain because a single large update can overly increase/decrease action probabilities.

18.3 Trust Region Policy Optimization (TRPO)

Policy gradient methods are susceptible to large updates that can destabilize training. If a policy update is too large, it can cause a catastrophic drop in performance.

TRPO addresses this by introducing the concept of a “trust region.” It aims to find the largest possible step that improves the policy, while ensuring that the new policy is not too far from the old policy. “Too far” is measured by the Kullback-Leibler (KL) divergence.

$$\begin{aligned} \max_{\theta} \mathbb{E}_{\pi_{\theta_{\text{old}}}} \left[\frac{\pi_\theta(A_t|S_t)}{\pi_{\theta_{\text{old}}}(A_t|S_t)} A_{\pi_{\theta_{\text{old}}}}(S_t, A_t) \right] \\ \text{s.t. } \mathbb{E}_{\pi_{\text{old}}} [D_{KL}(\pi_{\text{old}} || \pi_\theta)] \leq \delta \end{aligned}$$

- δ is a small hyperparameter controlling the maximum allowed KL divergence.

How do we ensure we don't go too far:

Using Hessian (second order derivative) TRPO ensures it doesn't violate the KL divergence constraint δ .

But! This method is complex to implement due to the second-order optimization required to handle the constraint.

18.4 Proximal Policy Optimization (PPO)

PPO is a popular, simpler alternative to TRPO that uses a clipped surrogate objective to achieve a similar effect of constraining policy updates.

$$J^{PPO}(\theta) = \max_{\theta} \mathbb{E}_{\pi_{\theta_{old}}} \left[\min \left(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right]$$

- ϵ is a hyperparameter (e.g., 0.1 or 0.2) that defines the “clipping range.”
- If Advantage is positive:

$$J^{PPO}(\theta) = \mathbb{E}_{\pi_{\theta_{old}}} \left[\min \left(\frac{\pi_{\theta}(A_t|S_t)}{\pi_{\theta_{old}}(A_t|S_t)}, 1 + \epsilon \right) A_{\pi_{\theta_{old}}}(S_t, A_t) \right]$$

- If Advantage is negative:

$$J^{PPO}(\theta) = \mathbb{E}_{\pi_{\theta_{old}}} \left[\max \left(\frac{\pi_{\theta}(A_t|S_t)}{\pi_{\theta_{old}}(A_t|S_t)}, 1 - \epsilon \right) A_{\pi_{\theta_{old}}}(S_t, A_t) \right]$$

where $r_t(\theta) = \frac{\pi_{\theta}(A_t|S_t)}{\pi_{\theta_{old}}(A_t|S_t)}$ is the probability ratio and \hat{A}_t is the advantage estimate. The ‘clip’ function prevents the policy from making excessively large changes, thereby avoiding catastrophic updates.

- If $A_t > 0$ (good action): We want to increase the probability of this action. The objective takes the minimum of $r_t(\theta)A_t$ and $(1 + \epsilon)A_t$. If $r_t(\theta)$ goes above $1 + \epsilon$, it gets clipped. This prevents updates from becoming too large and making the new policy vastly different from the old one, even if the action is very good.
- If $A_t < 0$ (bad action): We want to decrease the probability of this action. The objective takes the minimum of $r_t(\theta)A_t$ and $(1 - \epsilon)A_t$. If $r_t(\theta)$ goes below $1 - \epsilon$, it gets clipped. This prevents the policy from being penalized too much for actions that were just slightly worse than average. It avoids pushing the policy too far away from the old policy.

18.4.1 PPO Algorithm

The PPO algorithm is similar to Batch A2C but uses the PPO clipped objective for the actor update.

1. Sample trajectories $\{\tau_i\}$ from the current policy $\pi_{\theta_{old}}$.
2. Update the critic network (as in A2C).
3. Evaluate the advantage function (as in A2C).
4. Update the actor network for multiple epochs using the PPO objective on the same batch of data.
5. Set the old policy parameters: $\theta_{old} \leftarrow \theta$.

19 Off-Policy Actor-Critic Methods

19.1 Online Actor-Critic (On-Policy)

The basic actor-critic algorithm updates after each single time step. While simple, it has issues with instability when used with deep neural networks.

1. Take action $A_t \sim \pi_\theta(\cdot|S_t)$, observe $S_t, A_t, R_{t+1}, S_{t+1}$.
2. **Critic Update:** Update the value function \hat{v}_w towards the target $y_t = R_{t+1} + \gamma\hat{v}_w(S_{t+1})$.
3. **Actor Update:** Update the policy π_θ using the TD-error as the advantage signal: $\hat{A}(S_t, A_t) = y_t - \hat{v}_w(S_t)$.

In the on-policy version, the data used for learning comes from the same policy that's being updated. This can be inefficient because each data sample can only be used once.

19.2 Why Off-Policy Actor-Critic?

- Off-policy allows you to learn from any past experience, not just from the current policy's behavior.
- This is made possible through the use of a replay buffer (which stores past experiences) and target networks (which stabilize learning).
- Instead of learning the state value $V(s)$, which is tricky off-policy, we learn the action value $Q(s, a)$, as Q can be estimated from replay buffer samples without bias from old policies.

19.3 Off-Policy Actor-Critic

This is a more modern approach that combines the benefits of actor-critic methods with a replay buffer, allowing it to be off-policy. To support off-policy learning, it typically uses a Q-function instead of a V-function.

1. Take action A_t from the current policy π_θ and store the transition $(S_t, A_t, R_{t+1}, S_{t+1})$ in a replay buffer \mathcal{B} .
2. Sample a mini-batch of transitions from \mathcal{B} .
3. **Critic Update:** Update the critic's Q-network \hat{q}_w towards a target value y_i .

$$y_i = R_{i+1} + \gamma\hat{q}_w^-(S_{i+1}, A_{i+1}^{\pi_\theta})$$

where $A_{i+1}^{\pi_\theta} \sim \pi_\theta(\cdot|S_{i+1})$ is a sampled action from the current policy.

4. **Actor Update:** Update the actor π_θ to maximize the Q-values.

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i \nabla_\theta \log \pi_\theta(A_i^{\pi_\theta}|S_i) \hat{q}_w(S_i, A_i^{\pi_\theta})$$

5. **Target Network Update:** The target critic network w^- is updated slowly using a soft update: $w^- \leftarrow \tau w^- + (1 - \tau)w$.

19.4 Deep Deterministic Policy Gradient (DDPG)

As discussed in Section 16.4 DDPG is an off-policy, actor-critic algorithm for continuous action spaces. It uses a deterministic policy $\mu_\theta(s)$ for acting, with added noise for exploration, $A_t = \mu_\theta(S_t) + \mathcal{N}_t$.

19.4.1 Pseudocode

1. Initialize actor and critic networks μ and Q with parameters θ and w .
2. Initialize target networks μ^- and Q^- with parameters $\theta^- \leftarrow \theta$ and $w^- \leftarrow w$.
3. Initialize replay buffer \mathcal{B} .
4. For each episode:
 - (a) Initialize a random process \mathcal{N} for exploration.
 - (b) Receive initial state S .
 - (c) For $t = 1, T$:
 - i. Select action $A = \mu_\theta(S) + \mathcal{N}_t$.
 - ii. Execute action A and observe new state S' and reward R .
 - iii. Store transition (S, A, R, S') in \mathcal{B} .
 - iv. Sample a random minibatch of N transitions from \mathcal{B} .
 - v. For each sampled transition (S_i, A_i, R_i, S'_i) , calculate the target value:

$$y_i = R_i + \gamma Q_{w^-}(S'_i, \mu_{\theta^-}(S'_i))$$

- vi. Update the critic network by minimizing the loss:

$$L(w) = \frac{1}{N} \sum_i (y_i - Q_w(S_i, A_i))^2$$

- vii. Update the actor network using the sampled policy gradient:

$$\nabla_\theta J \approx \frac{1}{N} \sum_i \nabla_A Q_w(S_i, A)|_{A=\mu_\theta(S_i)} \nabla_\theta \mu_\theta(S_i)$$

Now the objective is to maximize the critic.

- viii. Soft update the target networks:

$$w^- \leftarrow \tau w + (1 - \tau)w^-$$

$$\theta^- \leftarrow \tau \theta + (1 - \tau)\theta^-$$

1. Broken Version — Off-Policy Actor-Critic v0.5

What it does:

- Uses a replay buffer but still tries to learn $V(s)$ directly from old (off-policy) samples.
- Critic target:

$$y_t = r_{t+1} + \gamma \hat{v}_{w^-}(s_{t+1})$$

- Actor uses advantages computed from V estimates.

Why it's broken:

- Off-policy issue: You can't directly learn $V(s)$ from replay buffer samples collected by older policies.
- The state and action distributions no longer match the current policy, making Bellman targets biased.
- This leads to incorrect value targets and unstable learning.

2. The Fix — Learn $Q(s, a)$ Instead (Off-Policy Actor-Critic v0.8)

Key change:

- Replace V -learning with a Q -learning style target:

$$y_t = r_{t+1} + \gamma \hat{q}_w(s_{t+1}, a_{t+1} \sim \pi_\theta(\cdot|s_{t+1}))$$

- Here, a_{t+1} is sampled from the *current* policy, not the replay buffer.

Why it works:

- $Q(s, a)$ Bellman targets remain valid off-policy — you only need the next action from the current policy.
- Advantage can now be computed as:

$$\hat{A}(s_i, a_i) = \hat{q}_w(s_i, a_i) - \hat{v}(s_i)$$

- With:

$$\hat{v}(s_i) = \sum_a \pi_\theta(a|s_i) \hat{q}_w(s_i, a)$$

20 Soft Actor-Critic (SAC)

Standard RL: “I’ll keep doing what seems best right now.”

SAC: “I’ll do what seems best on average, but keep some healthy randomness so I don’t miss a better long-term strategy.”

20.1 Maximum Entropy Reinforcement Learning

Instead of just maximizing reward, maximum entropy reinforcement learning adds an entropy bonus to the objective:

$$J(\pi) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t (R_{t+1} + \alpha \mathcal{H}(\pi(\cdot|S_t))) \right]$$

Here, the entropy of the policy π at state S_t is defined as:

$$\mathcal{H}(\pi(\cdot|S_t)) = - \sum_a \pi(a|S_t) \log \pi(a|S_t)$$

Entropy measures the randomness in the policy’s action choices. The parameter α is the temperature, which determines how much weight is given to the entropy bonus relative to the reward.

- **SAC is a model-free, off-policy, actor-critic algorithm that aims to find a policy that not only maximizes the reward but also has high entropy.**
- This encourages exploration and leads to more stable policies.
- The objective function is augmented with an entropy term:

$$\pi^* = \arg \max_{\pi} \sum_t \mathbb{E}_{\rho_{\pi}} [r(s_t, a_t) + \alpha \mathcal{H}(\pi(\cdot|s_t))]$$

where $\mathcal{H}(\pi(\cdot|s_t)) = \mathbb{E}_{a_t \sim \pi} [-\log \pi(a_t|s_t)]$ is the policy’s entropy and,

- α is the temperature parameter that controls the trade-off between reward and entropy.

20.2 Soft Value Functions

The Bellman equations are modified to include the entropy term, resulting in “soft” value functions.

Generally :

$$V^\pi(s_t) = \mathbb{E}_{a_t \sim \pi}[Q^\pi(s_t, a_t)]$$

Soft state-value:

$$V^\pi(s_t) = \mathbb{E}_{a_t \sim \pi}[Q^\pi(s_t, a_t) - \alpha \log \pi(a_t|s_t)]$$

Soft action-value:

$$Q^\pi(s_t, a_t) = r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1}}[V^\pi(s_{t+1})]$$

- In optimization, a **soft maximum** (or *soft value*) means we do not pick the single best action deterministically.
- Instead, we average over all actions, weighting them by their probability under π , and encourage a broader distribution.
- This avoids hard-max selection from Q-values (as in DQN), which can lead to premature convergence.

20.3 SAC Algorithm

SAC uses a stochastic policy and two Q-networks to mitigate overestimation bias. It is particularly effective for continuous control tasks.

1. Sample action A_t from the policy π_θ , execute it in the environment, and store the transition in a replay buffer D .
2. Sample a mini-batch of transitions from D .
3. **Critic Update:** Update both Q-networks to minimize the MSE loss against a soft Q-target, which is constructed using the minimum of the two Q-networks to reduce overestimation.

$$J_Q(\phi) = \mathbb{E}_{(s,a,s') \sim D} \left[\left(Q_\phi(s, a) - \left(r + \gamma \left(\min_{j=1,2} Q_{\phi_j^-}(s', \tilde{a}') - \alpha \log \pi_\theta(\tilde{a}'|s') \right) \right) \right)^2 \right]$$

where \tilde{a}' is a sampled action from the current policy.

4. **Actor Update:** Update the actor to minimize the expected Kullback-Leibler (KL) divergence the softmax over Q-Values and the policy π_θ to maximize the objective $J_\pi(\theta)$.

$$J_\pi(\theta) = \mathbb{E}_{s \sim D, a \sim \pi_\theta}[Q_\phi(s, a) - \alpha \log \pi_\theta(a|s)]$$

For continuous actions, a **reparameterization trick** is used to make this differentiable.

5. Update the temperature parameter α (optional) and perform soft updates on the target networks.

21 Model-Based Reinforcement Learning

21.1 Overview

Model-based RL is a paradigm where an agent learns a model of the environment’s dynamics from its experience. This model is then used for planning to inform and construct better policies or value functions. This contrasts with model-free RL, which learns policies and value functions directly from experience without explicitly learning a model.

21.1.1 Model-Free vs. Model-Based

- **Model-Free RL:** The agent learns a policy $\pi(a|s)$ or a value function $Q(s, a)$ directly from interactions with the environment. Examples include DQN and Policy Gradient methods.
- **Model-Based RL:** The agent learns a dynamics model $p(s', r|s, a)$ that predicts the next state s' and reward r given the current state s and action a . This model is then used to simulate new experiences and improve the policy.

21.2 Advantages and Disadvantages

21.2.1 Advantages

- **High Sample Efficiency:** By simulating experience from the learned model, the agent can learn a good policy with far fewer real-world interactions.
- **Rapid Adaptation:** The model can be quickly updated with new data, allowing the agent to adapt to changes in the environment or new tasks.
- **Efficient Model Learning:** Learning a predictive model is a supervised learning problem, which is often more stable and data-efficient than learning a policy via trial and error.

21.2.2 Disadvantages

- **Two Sources of Error:** Performance is sensitive to both the accuracy of the learned model and the effectiveness of the planning algorithm used with the model. Errors can compound over long planning horizons.

22 Model Learning

Learning a dynamics model is a supervised learning task. Given a dataset of transitions $\mathcal{D} = \{(s_i, a_i, r_{i+1}, s'_{i+1})\}$, the model learns to predict the next state and reward.

- **Reward Model:** A regression model predicts the reward, $r_{i+1} = f_{r\eta}(s_i, a_i)$.
- **State Transition Model:** A regression or density estimation model predicts the next state, $s'_{i+1} \sim f_{s\eta}(s_i, a_i)$.

22.1 Types of Dynamics Models

- **Deterministic Models:** The model outputs a single predicted next state, $s'_{i+1} = f_{s\eta}(s_i, a_i)$. This is suitable for environments with low stochasticity.
- **Stochastic Models:** The model outputs a probability distribution over the next state, from which a sample can be drawn. This is more robust for uncertain environments.

22.2 Combining Model-Free RL with a Model - Background Planning

A common approach is to use a learned model to generate synthetic data, which is then used to train a separate model-free RL algorithm. This is a powerful way to boost sample efficiency. However, major challenges is:

- **Distribution shift:** the model's accuracy is limited to the states it has seen.
- **Long synthetic rollouts** can accumulate prediction errors, leading to a "bad" policy.

22.2.1 Model-Based RL: Broken \rightarrow Fixed Progression (MBPO)

1. Naïve Model-Based RL (v0.5) — Broken

Algorithm:

1. Collect initial dataset \mathcal{D}_{env} in the real environment using policy π_{θ} .
2. **Loop:**
 - (a) **Model learning:** Fit dynamics model $\hat{P}(s'|s, a)$ on \mathcal{D}_{env} .
 - (b) **Policy updates:**
 - i. Simulate full trajectories entirely in \hat{P} starting from the initial state distribution.
 - ii. Update π_{θ} to maximize returns in \hat{P} .
 - (c) **No new real data collected here — major flaw.**

Why broken:

- Severe **distribution shift** and **compounding model errors** when rollouts go far from seen states.

2. Interleaving with Real Data (v0.8)

Algorithm:

1. Initialize \mathcal{D}_{env} with real transitions from an initial policy π_{θ} .
2. **Loop:**
 - (a) **Model learning:** Train dynamics model \hat{P} on \mathcal{D}_{env} .
 - (b) **Policy updates:**
 - i. Simulate trajectories entirely in \hat{P} to create $\mathcal{D}_{\text{model}}$.
 - ii. Update π_{θ} using $\mathcal{D}_{\text{model}}$.
 - (c) **Collect real data:** Deploy π_{θ} in the real environment for N steps and add to \mathcal{D}_{env} . ← new change

Improvement:

- Model continuously updated with states from current policy.

Still broken:

- Long simulated rollouts still accumulate large prediction errors.

3. Short Model Rollouts from Real States (v1.0)

Algorithm:

1. Initialize \mathcal{D}_{env} with real transitions from an initial policy π_{θ} .
2. **Loop:**
 - (a) **Collect real data:** Deploy π_{θ} in the real environment for N steps and add to \mathcal{D}_{env} .
 - (b) **Model learning:** Train dynamics model \hat{P} on \mathcal{D}_{env} .
 - (c) **Policy updates:**

- i. Sample real states from \mathcal{D}_{env} .
- ii. Perform **short rollouts** (1–5 steps) in \hat{P} from those states to create $\mathcal{D}_{\text{model}}$. ← new change
- iii. Train π_θ on $\mathcal{D}_{\text{model}} \cup \mathcal{D}_{\text{env}}$. ← new change

Improvement:

- Reduces distribution shift since all rollouts start from real states.
- Limits compounding errors with short horizons.

22.2.2 Dyna-Q Style Algorithm - MBPO

This algorithm uses a model to supplement real-world experience, often in a continuous loop.

Algorithm:

1. Run base policy $\pi_0(a|s)$ in the real environment to collect $D_{\text{env}} = \{(s, a, r, s')\}_i$.
2. Learn environment model $f_\eta(s, a)$ using D_{env} .
3. Pick states S_i from D_{env} , use $f_\eta(s, a)$ to make short rollouts; add to D_{model} .
4. Use both D_{env} and D_{model} to update $\pi_\theta(a|s)$ via any off-policy RL algorithm (e.g., DQN, SAC). ← picks states from both env and the model
5. Deploy $\pi_\theta(a|s)$ in the real environment, appending (s, a, r, s') to D_{env} and repeat.

This approach uses short rollouts from the model to avoid accumulating errors and combines real data to anchor learning, leading to more stable performance.

Improvement:

- Combines sample efficiency of model rollouts with realism of real data.

State coverage advantage:

- Model rollouts may reach states not yet visited in the real environment (from s_0).
- Expands exploration coverage without costly real-environment interaction.

22.2.3 Direct Policy Learning with a Model

This approach directly optimizes the policy by backpropagating gradients through the learned dynamics model. The key idea is to treat the environment model as part of the computational graph.

1. Run the policy π_θ in the real environment to collect a dataset \mathcal{D}_{env} .
2. Learn a dynamics model $f_\eta(s, a)$ from \mathcal{D}_{env} .
3. To improve the policy, unroll the policy's actions through the model's dynamics for a few steps. The reward from this unrolling is used to compute a loss.
4. Backpropagate the gradients of this loss through the dynamics model and into the policy network to update the policy parameters θ .
5. Repeat the process by running the updated policy in the real environment to collect new data.

This method can be very sample-efficient, but it can suffer from issues like vanishing or exploding gradients when unrolling over many time steps, and its effectiveness is highly dependent on the accuracy of the learned model.

Benefit:

- End-to-end differentiable optimization.

Downside:

- Sensitive to model inaccuracies.
- Gradient instability for large H (Vanishing/Exploding gradients).

Summary Table

Version	Change	Main Benefit	Remaining Problem
v0.5	Train policy only in model	High sample efficiency	Severe shift, compounding error
v0.8	Periodic real data collection	Reduces shift	Still long rollout errors
v1.0	Short rollouts from real states	Limits error growth	Short horizon only
MBPO	Mix real + model data	Stable + efficient	Needs tuning
Direct BP	Backprop through model	Gradient-based planning	Numerical instability

23 Model Based - Decision-Time Planning

The goal is to find the action sequence that maximizes the sum of future reward.

In contrast to background planning, which improves the policy over time, decision-time planning uses the model to select the best action for the current state. The model is essentially used as a simulator to find an optimal action sequence. The goal is to find an action sequence that maximizes the sum of future rewards:

$$a_1, \dots, a_T = \arg \max_{a_1, \dots, a_T} \sum_{t=1}^T \gamma^{t-1} r_\eta(s_t, a_t) \quad \text{s.t.} \quad s_{t+1} = f_\eta(s_t, a_t)$$

23.1 Stochastic Optimization

- **Random Shooting Method:** This method samples a large number of random action sequences and evaluates each one using the learned model. The first action of the best-performing sequence is then executed. This is simple but can be inefficient in high-dimensional or long-horizon problems.
- **Cross-Entropy Method (CEM):** CEM is an iterative optimization algorithm that refines a distribution of action sequences.
 1. Sample a number of action sequences from a distribution (e.g., a Gaussian).
 2. Evaluate the performance of each sequence using the model.
 3. Select the top-performing “elite” sequences.
 4. Fit a new distribution to these elite sequences to generate the next batch of samples.
 5. Repeat until convergence.

23.2 Monte Carlo Tree Search (MCTS)

MCTS is a planning algorithm that explores a search tree to find the best action. It is commonly used in game-playing (e.g., AlphaGo).

1. **Selection:** Starting from the root, traverse the tree by choosing actions that balance exploration (visiting less-explored nodes) and exploitation (choosing nodes with high value estimates) until a leaf node s_L is reached. A common selection strategy is the UCT formula:

$$Q(s, a) + c\sqrt{\frac{\log N(s)}{N(s, a)}}$$

2. **Expansion:** If the selected leaf node has not been fully expanded, create a new child node for an unvisited action.
3. **Simulation:** From the new leaf node, run a simulation (a “rollout”) using a fast default policy until a terminal state is reached.
4. **Backpropagation:** The result of the simulation is used to update the value estimates and visit counts of all nodes from the leaf back to the root.

23.3 Model-Based RL with Decision-Time Planning

This approach, often called Model Predictive Control (MPC), combines a learned model with decision-time planning.

1. Collect data using an initial policy and learn a dynamics model $f_\eta(s, a)$.
2. For a given state, use the learned model and a planning algorithm (e.g., CEM or MCTS) to find the best sequence of actions.
3. **Execute only the first action of this sequence in the real environment.**
4. Observe the next state, add the transition to the dataset, and repeat the process for the new state.

This is a robust method because it uses the model for planning, but only commits to the first action, allowing it to replan at every step with new, real-world information.

23.4 Uncertainty in RL

23.4.1 Sources

- **Epistemic (model uncertainty):** Due to limited data; reduces with exploration.
- **Aleatoric (intrinsic randomness):** Due to stochasticity in environment; cannot be reduced.

23.4.2 In Model-Based RL

Learned dynamics model $f_\eta(s, a)$ may be inaccurate in unexplored regions → planning might exploit these errors.

- **Example:** MPC planning might choose actions leading into regions of high prediction uncertainty (red circled areas in plots).

23.4.3 Mitigation

- Use ensembles or Bayesian models to estimate epistemic uncertainty.
- Penalize reward in high-uncertainty states (uncertainty-aware planning).
- Combine real and simulated rollouts cautiously; trust short rollouts more in uncertain regions.

23.5 Uncertainty in Neural Networks

23.5.1 Estimate model uncertainty (how unsure the parameters are)

How probable is θ given the dataset. This probability will have an entropy, which will say how uncertain the model is.

- Find $p(\theta \mid \mathcal{D})$ — distribution of parameters given data.
- Entropy of $p(\theta \mid \mathcal{D})$ = how uncertain the model is.

Prediction Average over possible parameters:

$$p(s_{t+1} \mid s_t, a_t, \mathcal{D}) = \int p(s_{t+1} \mid s_t, a_t, \theta) p(\theta \mid \mathcal{D}) d\theta$$

Avoid overconfident exploitation of model errors; uncertainty estimates guide safer, more reliable exploration and planning.

23.5.2 Bootstrap Ensembles

If you have K environment models of the environment and they all predict the same state transition - $p(s_{t+1} \mid s_t, a_t, \theta)$, the uncertainty is low; if their predictions differ, the uncertainty is high.

Method:

- **Dataset:** $\mathcal{D} = \{(s, a, r, s')\}_i$.
- Train K models.
- For each model, sample batches $\mathcal{D}_1, \mathcal{D}_2, \dots, \mathcal{D}_K$ with replacement from \mathcal{D} .

Before: Plan using one model:

$$a_{1:T} = \arg \max_{a_{1:T}} \sum_{t=1}^T \gamma^t r_{\eta}(s_t, a_t) \quad \text{s.t.} \quad s_{t+1} = f_{\eta}(s_t, a_t)$$

Now: Plan using K models (ensemble) and average rewards:

$$a_{1:T} = \arg \max_{a_{1:T}} \frac{1}{K} \sum_{i=1}^K \sum_{t=1}^T \gamma^t r_{\eta_i}(s_{t,i}, a_t)$$

where $s_{t+1,i} = f_{\eta_i}(s_{t,i}, a_t)$.

Key idea: Choose actions that perform well on average across all models \rightarrow robust to model uncertainty.

Diversity in environment model comes from:

- Different neural network initializations.
- Different training samples for each model.

Benefit: Differences between models give an estimate of epistemic uncertainty.

23.5.3 Bayesian Neural Networks (BNNs)

- **Goal:** Capture model uncertainty by learning a *distribution* over network weights.
- **Parameter Uncertainty:** Estimate the posterior $p(\theta \mid \mathcal{D})$.
- **Uncertainty Measure:** Entropy of $p(\theta \mid \mathcal{D})$ indicates model uncertainty.
- **Prediction:** Average over possible parameter values:

$$p(s_{t+1} \mid s_t, a_t, \mathcal{D}) = \int p(s_{t+1} \mid s_t, a_t, \theta), p(\theta \mid \mathcal{D}), d\theta$$

- **Benefit:** Earlier we had scalar values for the parameter (θ), now we will have μ and σ from which we can see for a given input how much the output can vary. Here we are directly encoding the probability of θ .

23.6 Latent Space Models

In RL, especially when observations are high-dimensional (e.g., images), we often compress them into a low-dimensional latent state space s_t using an encoder.

Then, the environment's components are modeled separately in this latent space:

- **Observation Model** $p(o_t \mid s_t)$: Maps latent state s_t back to observation o_t . Handles high-dimensional data (like pixels). Doesn't need to model temporal dynamics.
- **Dynamics Model** $p(s_{t+1} \mid s_t, a_t)$: Predicts the next latent state from the current one and the action. Low-dimensional but must capture temporal dependencies.
- **Reward Model** $p(r_t \mid s_t, a_t)$: Predicts reward from latent state and action.

24 Safe Reinforcement Learning

24.1 Overview of Safety Problems

Safe Reinforcement Learning (Safe RL) is a field dedicated to ensuring that an agent's behavior adheres to constraints and avoids undesirable, potentially harmful, outcomes. Standard RL aims to maximize a single reward function, which may not be sufficient to guarantee safety. Key issues in Safe RL include:

- **Optimal but Unsafe Behavior:** An optimal policy can lead to high rewards but also high-risk actions that should be avoided.
- **Safe Exploration:** The agent must act safely not only at convergence but also during the training process, where it actively explores the environment.
- **Irreversible States:** Preventing the agent from entering states from which a return to a safe state is impossible.

24.2 Constrained Markov Decision Processes (CMDP)

To formalize safety, Safe RL often uses the Constrained Markov Decision Process (CMDP) framework. A CMDP extends a standard MDP by adding one or more cost functions and associated thresholds.

- **Classic MDP:** $\langle S, A, P, R, \gamma \rangle$, where the objective is to maximize $V(s_0) = \mathbb{E}_\pi [\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)]$.
- **CMDP:** A CMDP adds a cost function $C : S \times A \rightarrow \mathbb{R}$ and a set of thresholds b_i . The objective is to maximize the expected return subject to a set of constraints on the expected cumulative cost.

24.2.1 Machiavelli Benchmark

- The Machiavelli Benchmark is a safety-focused evaluation framework for reinforcement learning and language model agents, designed to test whether they can achieve goals without engaging in unethical or harmful behaviors.
- **Observation:** narrative state (NPC offers collusion under threat).
- **Actions:**
 - *Deceptive collusion* (advance ambitions quickly).
 - *Honest cooperation* (pretend to help, secretly oppose).
 - *Self-sacrifice* (enter the mines to investigate).

24.2.2 Optimally Safe Policies

An optimally safe policy π_c^* is one that achieves the highest possible return while satisfying all safety constraints.

$$\pi_c^* = \arg \max_{\pi_c \in \Pi_c} V(s_0)$$

where Π_c is the set of all policies that satisfy the constraints. These constraints can be on the expected total cost or on the probability of a catastrophic event.

- **Expected Safety:** The expected cumulative cost must be below a threshold:

$$f_c(\pi) = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t C(s_t, a_t) \right] \leq b_i$$

- **Almost Surely Safe:** The probability of the cumulative cost exceeding the threshold is zero.

$$P_\pi \left(\sum_{t=0}^{\infty} \gamma^t C(s_t, a_t) \leq b_i \right) = 1$$

24.3 Solving CMDPs: Methods

A variety of methods exist to solve CMDPs, including policy constraint algorithms, reward shaping, and modified action selection strategies.

24.3.1 Constrained Policy Optimization (CPO)

CPO is an extension of the TRPO algorithm that adds a safety constraint. It finds a policy update that maximizes the objective while ensuring the KL divergence from the old policy and the cumulative cost are within specified bounds.

$$\max_{\theta} \mathcal{J}^{TRPO}(\theta) = \mathbb{E}_{\pi_{old}} \left[\frac{\pi_{\theta}(a_i | s_i)}{\pi_{old}(a_i | s_i)} A_{\pi_{old}}(s_i, a_i) \right]$$

$$\text{s.t. } \mathbb{E}_{\pi_{\theta_{\text{old}}}}[D_{KL}(\pi_{\theta_{\text{old}}}||\pi_{\theta})] \leq \delta \quad \text{and} \quad \mathbb{E}_{\pi_{\theta_{\text{old}}}}[C(s, a)] \leq b_i$$

Insights from TRPO:

- Monotonic Performance Improvement (each time the policy is updated, its expected return (or value) will be at least as good as, and typically better than, the previous policy).
- Satisfy constraints

24.3.2 Primal-Dual Methods (Lagrangian)

Lagrangian methods convert a constrained optimization problem into an unconstrained one by introducing a Lagrange multiplier λ for each constraint. The objective is to find a saddle point of the Lagrangian function.

$$\min_{\lambda} \max_{\theta} \mathcal{L}(\theta, \lambda) = \mathbb{E}_{\pi}[R(s, a)] - \lambda (\mathbb{E}_{\pi}[C(s, a)] - b)$$

The method iteratively updates the policy to maximize this objective and the multiplier to enforce the constraint.

1. Initialize policy parameters θ and Lagrange multiplier $\lambda \geq 0$.
2. For each iteration, compute policy and value losses. The actor loss is augmented with the cost term:

$$\mathcal{L}_{\theta} = -\mathbb{E}_{\pi}[\log \pi_{\theta}(a|s) \cdot A] - \lambda (\mathbb{E}_{\pi}[C(s, a)] - b)$$

3. Update θ via a gradient step on \mathcal{L}_{θ} .
4. Update the Lagrange multiplier: $\lambda \leftarrow \lambda + \alpha (\mathbb{E}_{\pi}[C(s, a)] - b)$. This increases λ if the cost constraint is violated, penalizing future cost increases.

Primal-Dual Methods: Example

- **Primal Problem:** The original constrained optimization problem is

$$\max f(x) = 3x \quad \text{subject to } x - 4 \leq 0 \quad (\text{or } x \leq 4).$$

The optimal value is achieved at $x = 4$, giving $f(4) = 12$.

- **Unconstrained Lagrangian:** The constraint is incorporated using a Lagrange multiplier λ , creating the Lagrangian:

$$L(x, \lambda) = 3x - \lambda(x - 4).$$

- **Dual Problem:** The dual method reverses the order of optimization:

$$\min_{\lambda} \max_x [3x - \lambda(x - 4)].$$

We first maximize with respect to x for fixed λ , and then minimize with respect to λ .

Solving via Table Analysis

The values of $L(x, \lambda)$ are computed in a table for $x = 3, 4, 5$ and $\lambda = 2, 3, 4$:

λ	$x = 3$	$x = 4$	$x = 5$
2	10	12	13
3	12	12	12
4	13	12	11

Optimization Steps:

- For each λ , find the maximum value in the row (max over x).
- For $\lambda = 2$: $\max\{10, 12, 13\} = 13$
- For $\lambda = 3$: $\max\{12, 12, 12\} = 12$
- For $\lambda = 4$: $\max\{13, 12, 11\} = 13$
- Now, minimize these maxima over λ : $\min\{13, 12, 13\} = 12$

So, $\lambda^* = 3$ yields the minimum of the maximum Lagrangian value.

Conclusion

- $\lambda = 3$ is the optimal solution to the dual problem,
- The optimal dual value equals the optimal primal value (12).
- This confirms strong duality for this simple linear program.

Summary Table:

λ	$\max_x L(x, \lambda)$
2	13
3	12
4	13

Optimal $\lambda^* = 3$, value = 12.

Limitations of Lagrangian:

You can't guarantee that the safety constraint will be satisfied. For example, your budget may be to not exceed 5 speed violations in an autonomous car, but with a Lagrangian method there can still be a 6th violation, because it only enforces the limit on average rather than in every episode.

24.3.3 State Augmentation (SauteRL)

This method reformulates a CMDP as a standard MDP by augmenting the state with a “safety budget.” The budget decreases with each cost incurred.

- **Augmented State:** The state becomes $\tilde{S} = S \times Z$, where Z is a variable that tracks the remaining safety budget.
- **Modified Reward:** The agent receives the original reward $R(s, a)$ only if the safety budget is non-negative ($z \geq 0$). Otherwise, the reward is zero.
- **Budget Update:** The budget is updated at each step: $z_{t+1} = z_t - f_c(s_t, a_t)$, with an initial budget $z_0 = b_i$.

A standard RL algorithm (e.g., PPO or SAC) can then be used to solve this augmented MDP. This allows the policy to generalize across different constraints by simply changing the initial budget.

24.3.4 Action Selection (MASE)

The Meta-Algorithm for Safe Exploration (MASE) is a high-level approach that guides exploration to prevent unsafe actions. It assumes a safety margin and an emergency stop action.

- **Uncertainty Quantification:** A model estimates the cost of an action and its uncertainty, $\Gamma(s, a)$.
- **Safe Actions:** A set of “safe” actions is defined as those where the estimated cost plus uncertainty does not exceed the safety threshold, i.e., $A^+ = \{a \mid \mu(s, a) + \Gamma(s, a) \leq b_i\}$.
- **Exploration:** The agent’s policy is restricted to choosing actions from the safe set A^+ . If no safe actions are available, an emergency action (e.g., stopping) is taken.

This ensures the agent acts safely during training while learning a policy that is also safe at convergence.

Definition of a Safe Action in MASE

An action is safe if:

$$a \in \mathcal{A} \quad \text{and} \quad \mu(s_i, a) + \Gamma(s_i, a) \leq b_i$$

where:

- $\mu(s_i, a) \rightarrow$ predicted safety cost (learned with a neural net, for example)
- $\Gamma(s_i, a) \rightarrow$ uncertainty bound on the prediction. (Γ - Capital Gamma)
- $b_i \rightarrow$ maximum allowed cost

Safety means: Even in the worst case (prediction + uncertainty), the cost stays within the limit.

How MASE Works (Step-by-Step)

At each time step h :

1. The agent is in state s_h .
2. It looks at the set of safe actions:

$$\mathcal{A}_h^+ = \{a \mid \mu(s_h, a) + \Gamma(s_h, a) \leq b\}$$

3. If safe actions exist ($\mathcal{A}_h^+ \neq \emptyset$):
 - Pick one (via your RL policy) and execute it.
 - Continue exploring.
4. If no safe action exists ($\mathcal{A}_h^+ = \emptyset$):
 - Trigger emergency stop action \hat{a} (from Assumption 2).
 - This keeps the agent in a safe state but comes with a large penalty so it learns to avoid such dead ends.

Why it works:

This guarantees almost sure safety because at every step:

- You either take a provably safe action.
- Or you stop immediately before violating constraints.

24.4 Frontiers in Safe RL

Safe RL is a rapidly evolving field with ongoing research in several key areas:

- **Safe Offline RL:** Learning safe policies from a fixed dataset without further interactions.
- **Safe Model-Based RL:** Using a learned model to plan for safe, long-horizon behaviors.
- **Safe Exploration:** Developing methods to learn and explore safely in the real world.
- **Reward Hacking:** Preventing agents from finding unintended ways to get high rewards that violate human-defined safety principles.

25 Exploration in RL

25.1 Random Exploration

Random exploration is a fundamental strategy for balancing the trade-off between **exploration** (trying new actions to find better policies) and **exploitation** (choosing the best-known action).

- **ϵ -Greedy:** This is a simple but effective strategy where, with a probability of $1 - \epsilon$, the agent selects the action with the highest estimated value (Q -value). With probability ϵ , it chooses a random action. The value of ϵ is often decayed over time to shift the balance from exploration to exploitation as the agent learns more about the environment.
- **In Deep RL:** For continuous action spaces, exploration can be achieved by adding noise (e.g., Gaussian noise) to the agent’s deterministic policy output. For discrete actions, a softmax function over the Q -values or an entropy regularization term in the policy’s objective can be used to encourage randomness.

25.2 Novelty Seeking Exploration

Novelty-seeking methods provide an **intrinsic reward** to the agent for visiting new or less-explored states. This encourages the agent to actively seek out new information.

- **Count-Based Exploration:** A bonus reward is given for visiting a state. The bonus is inversely proportional to the number of times the state has been visited. For example, the bonus could be $B(N(s, a)) = 1/\sqrt{N(s, a)}$, where $N(s, a)$ is the visit count for a state-action pair.
- **Pseudo-Counts:** For environments with large or continuous state spaces, it’s not feasible to keep a simple count. Pseudo-count methods use a density model to approximate how many times a state has been visited. The bonus is based on the novelty of a state as measured by this model.
- **Prediction Errors:** The agent’s bonus reward is proportional to the error of a predictive model (e.g., a neural network) trying to predict the next state. High prediction error suggests a novel state, encouraging the agent to explore it further.

25.3 Posterior Sampling Methods

These methods, also known as Bayesian exploration, model the uncertainty in the value function or model of the environment.

- **Upper Confidence Bound (UCB):** An optimistic exploration strategy that adds a bonus to the Q-value based on the uncertainty of the action’s value. The bonus term, $c\sqrt{\ln t/N_t(a)}$, encourages the agent to choose actions with high estimated values and actions that have been tried less often ($N_t(a)$ is low).
- **Bootstrapped DQN:** An ensemble of K Q-networks is trained on bootstrapped samples from the replay buffer. At the start of each episode, one Q-network is randomly selected, and the agent follows its greedy policy. This method effectively samples from a posterior distribution over Q-functions.

26 Transfer Learning in RL

Transfer learning aims to leverage knowledge from a source task to accelerate learning in a new, target task.

26.1 Domain Adaptation

Domain adaptation addresses the challenge of transferring a policy from a source environment (e.g., a simulator) to a target environment (e.g., the real world) where the observations or dynamics may be different.

- **Observation Adaptation:** This involves learning a representation of the state that is invariant to the domain (e.g., ignoring visual differences between a simulated and real robot arm). Adversarial methods are a common approach, where a domain classifier is trained to distinguish between source and target observations, while an encoder is simultaneously trained to fool the classifier, thus producing a domain-invariant representation.
- **Dynamics Adaptation:** When the physics of the two environments differ, the dynamics model can be adapted. This can be done by using a penalty on transitions that are unlikely in the target environment.

26.2 Domain Randomization

Instead of trying to match the source domain to the target, domain randomization varies the source domain’s parameters (e.g., friction, lighting, textures) to such a degree that the learned policy becomes robust enough to generalize to the real-world target domain.

26.3 Multi-Task Transfer

This involves learning a single policy that can solve multiple tasks simultaneously or a policy that can quickly adapt to a new task.

- **Contextual Policies:** A policy takes both the state and a task-specific context as input: $\pi(a|s, \text{context})$. The context can be a one-hot vector representing the task ID or a continuous vector embedding of the task.
- **Goal-Conditioned Policies:** The policy is trained to reach any of a set of goals, which are provided as a part of the input. Hindsight Experience Replay (HER) is a technique that can make this more data-efficient by “relabeling” past trajectories with the goal that was actually achieved, even if the intended goal was not met.

27 Frontiers in RL

RL is a rapidly advancing field with many open challenges and new research directions.

27.1 Meta-Learning

Meta-learning, or “learning to learn,” aims to create agents that can quickly adapt to new tasks from a distribution of tasks.

- **Model-Agnostic Meta-Learning (MAML)**: An algorithm that learns a good set of initial policy parameters such that a policy can be adapted to a new task with only a few gradient steps.

27.2 Inverse Reinforcement Learning (IRL)

IRL addresses the **reward specification problem** by inferring a reward function from expert demonstrations. Instead of manually designing a reward, the agent learns what the expert’s goal is by observing their behavior.

27.3 Hierarchical RL

Hierarchical RL introduces a hierarchy of policies to solve complex, long-horizon problems. A high-level policy selects a “sub-goal” or “option,” which is then executed by a lower-level policy. This provides a temporal abstraction that can make learning more efficient.

27.4 Foundation Models for RL

Recent research explores using large pre-trained models (e.g., LLMs) to provide high-level planning, world models, or reward functions. The idea is to pre-train a model on a vast amount of data and use it as a foundation for decision-making.

27.5 Continual RL

Continual RL focuses on agents that operate in vast, non-stationary worlds without a clear start or end, where rewards are sparse, and the environment changes over time. Challenges include irreversibility, rich observations, and catastrophic forgetting.

28 Sample Question and Answer

Question 1: REINFORCE with Baseline

The policy gradient update with a baseline is given by $\theta_{t+1} = \theta_t + \alpha \nabla_{\theta} \log \pi_{\theta}(A_t|S_t)(G_t - b(S_t))$, where G_t is the Monte-Carlo return and $b(S_t)$ is a learned value function $\hat{v}(S_t, w)$.

1. Why is it mathematically valid to subtract a baseline $b(S_t)$ that only depends on the state? Explain why this does not introduce bias into the policy gradient estimate. (2 sentences)

Answer: Subtracting a state-dependent baseline $b(S_t)$ is valid because its expected value, when multiplied by the score function $\nabla_{\theta} \log \pi_{\theta}(a|s)$, is zero. This property ensures that the expectation of the overall gradient estimate remains unchanged, meaning the update is still an unbiased estimate of the true policy gradient.

2. **What is the primary practical benefit of using a baseline in policy gradient methods, and how does the term $(G_t - b(S_t))$ relate to the advantage function? (3 sentences)**

Answer: The primary benefit of using a baseline is to reduce the variance of the policy gradient estimates, which leads to more stable and faster learning. The term $(G_t - b(S_t))$ is a Monte-Carlo estimate of the advantage function $A(S_t, A_t)$. This is because G_t is an estimate of the Q-value $q_\pi(S_t, A_t)$ and the baseline $b(S_t)$ is an estimate of the state-value $v_\pi(S_t)$.

Question 2: Proximal Policy Optimization (PPO)

The PPO algorithm uses a clipped surrogate objective:

$$\max_{\theta} \mathbb{E}_{\pi_{\theta_{old}}} \left[\min \left(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right]$$

where $r_t(\theta) = \frac{\pi_{\theta}(A_t|S_t)}{\pi_{\theta_{old}}(A_t|S_t)}$.

1. **What is the fundamental problem that the clip function is designed to prevent?**

Answer: The clip function is designed to prevent the policy from making excessively large updates by discouraging the probability ratio $r_t(\theta)$ from moving too far from 1. This avoids catastrophic performance collapses that can occur with unconstrained on-policy updates, effectively creating a “trust region.”

2. **How does this objective offer a simpler and more practical alternative to its predecessor, TRPO?**

Answer: PPO’s clipped objective provides a simpler, first-order optimization method to constrain the policy update. This avoids the complex and computationally expensive second-order optimization (like conjugate gradient) required by TRPO to solve its KL-divergence constraint.

Question 3: Soft Actor-Critic (SAC)

The SAC objective function is augmented with an entropy term:

$$\pi^* = \arg \max_{\pi} \sum_t \mathbb{E}_{\rho^{\pi}} [r(s_t, a_t) + \alpha \mathcal{H}(\pi(\cdot|s_t))]$$

1. **What is the primary role of the temperature parameter α ?** **Answer:** The temperature parameter α controls the trade-off between maximizing the cumulative reward and maximizing the policy’s entropy.
2. **What are two key benefits of maximizing policy entropy in addition to the expected reward, and how is this reflected in the soft Q-function update?** **Answer:** Maximizing entropy encourages better exploration by preventing premature convergence to a suboptimal policy and leads to more robust, stable policies. This is reflected in the soft Q-function update by adding the entropy of the next state’s policy to the target value: $y_i = r_i + \gamma(\min_j Q_j(s', a') - \alpha \log \pi(a'|s'))$. This modification ensures that the value functions account for the future entropy rewards.

Question 4: Offline Reinforcement Learning (CQL)

Offline RL learns from a fixed dataset, which presents a major challenge known as distribution shift.

1. **Briefly explain why out-of-distribution (OOD) actions are problematic for standard off-policy algorithms like Q-learning in the offline setting. Answer:** Standard Q-learning can produce erroneously high Q-values for OOD actions not present in the data because the function approximator has no data to constrain its estimates for those actions. The policy then learns to exploit these errors, leading to poor performance when deployed in the real world.
2. **The Conservative Q-Learning (CQL) algorithm adds a specific regularizer to the Bellman error loss. What is the high-level goal of this regularizer, and what two things does it do to the Q-function? Answer:** The high-level goal of the CQL regularizer is to combat Q-value overestimation for OOD actions. It explicitly pushes down the Q-values for actions that are likely OOD (sampled from a distribution $\mu(a|s)$). Simultaneously, it pushes up the Q-values for actions that are actually present in the dataset.

Question 5: Dynamic Programming vs. Model-Free Control

Consider two fundamental algorithms: Value Iteration (a Dynamic Programming method) and Q-Learning (a Model-Free Control method). Both aim to find the optimal action-value function q^* .

1. **What is the single most important difference in the assumptions these two algorithms make about the environment? (2 sentences) Answer:** Value Iteration assumes a complete model of the environment is known, meaning it requires the full transition probabilities $p(s', r|s, a)$ and reward function $R(s, a)$. Q-Learning, being model-free, makes no such assumption and can learn directly from sampled experiences (s, a, r, s') without knowing the underlying dynamics.
2. **Given that Value Iteration can be more computationally efficient for small, known environments, why is Q-Learning often more practical for solving large-scale or real-world problems? (2 sentences) Answer:** In most real-world problems, the true dynamics of the environment are unknown and far too complex to be explicitly defined or stored. Q-Learning's ability to learn from direct interaction and experience makes it applicable to these complex scenarios where building an explicit model is infeasible.

Question 6: On-Policy vs. Off-Policy Learning

Sarsa is a classic on-policy TD control algorithm, while Q-Learning is its off-policy counterpart.

1. **Explain the difference in their update rules and how it relates to their on-policy vs. off-policy nature. (2 sentences)**
 - **SARSA update:** $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t)]$
 - **Q-Learning update:** $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha[R_{t+1} + \gamma \max_{a'} Q(S_{t+1}, a') - Q(S_t, A_t)]$

Answer: Sarsa is on-policy because it uses the action A_{t+1} actually taken by the current policy to form its target, learning the value of the policy it is currently following. Q-Learning is off-policy because it uses the greedy action $\max_{a'} Q(S_{t+1}, a')$ to form its target, learning the value of the optimal (greedy) policy, regardless of which exploratory action was actually taken.

2. **What is the main advantage of off-policy learning (like Q-Learning) over on-policy learning (like Sarsa) in terms of data efficiency? Why might an on-policy method sometimes be preferred? (3 sentences) Answer:** The main advantage of off-policy learning is data efficiency, as it can learn the optimal policy from data generated by any behavior policy, including old data stored in a replay buffer. On-policy methods, in contrast, must discard data after each policy

update. However, an on-policy method might be preferred for its stability, as it directly learns about the consequences (the feedback you get from the environment is of the current policy and not using a different policy like off-policy) of the policy it's actually executing, which can lead to smoother convergence, especially in the presence of function approximation.

Question 7: Model-Based Reinforcement Learning

Model-based RL approaches learn a model of the environment to aid in policy learning.

1. **What is the primary advantage and a key disadvantage of model-based RL? (2 sentences) Answer:** The primary advantage is sample efficiency, as the learned model can generate many simulated experiences, reducing the need for real-world interaction. A key disadvantage is that the performance is limited by the accuracy of the learned model; errors in the model can be exploited by the policy, leading to poor real-world performance.
2. **How does the “Dyna-Q” style algorithm, which combines model-based and model-free elements, attempt to mitigate this key disadvantage? (2 sentences) Answer:** Dyna-Q mitigates the problem of model error by continuously updating the model with real experience collected from the environment. It combines planning with model-generated data and learning from real data, which helps to correct for inaccuracies in the model and ground the policy in reality.

Question 8: Deep Deterministic Policy Gradient (DDPG)

DDPG is an off-policy algorithm designed for continuous control.

1. **Why is DDPG considered an “Actor-Critic” algorithm? Briefly describe the roles of the actor and the critic. (2 sentences) Answer:** DDPG is an Actor-Critic algorithm because it uses two separate networks: an actor that learns a deterministic policy ($\mu_\theta(s)$) to select actions, and a critic that learns a Q-function ($Q_w(s, a)$) to evaluate those actions. The critic guides the actor's learning by providing a gradient signal, telling the actor how to adjust its policy to select actions that lead to higher Q-values.
2. **Since the actor's policy is deterministic, how does DDPG ensure sufficient exploration of the state-action space during training? (1 sentence) Answer:** DDPG ensures exploration by adding noise (typically from a stochastic process like Ornstein-Uhlenbeck or simple Gaussian noise) to the actions selected by the deterministic actor during the training phase.

Question 9: DDPG vs. SAC: A Mathematical Comparison

1. Deep Deterministic Policy Gradient (DDPG) Equations

DDPG is an off-policy actor-critic algorithm that uses a deterministic policy and a single critic.

Target value:

$$y^{\text{DDPG}} = r + \gamma Q_{\theta^-}(s', \mu_{\phi^-}(s'))$$

Where:

- $Q_\theta \rightarrow$ critic (estimates $Q(s, a)$).
- $Q_{\theta^-} \rightarrow$ target critic network.
- $\mu_\phi \rightarrow$ deterministic actor network.

- $\mu_{\phi^-} \rightarrow$ target actor network.
- $r \rightarrow$ reward received.
- $\gamma \rightarrow$ discount factor.

Critic loss:

$$L(\theta) = \mathbb{E}_{(s,a,r,s') \sim \mathcal{D}} \left[(Q_{\theta}(s, a) - y^{\text{DDPG}})^2 \right]$$

Actor loss:

$$J_{\mu}(\phi) = -\mathbb{E}_{s \sim \mathcal{D}} [Q_{\theta}(s, \mu_{\phi}(s))]$$

The actor maximizes the critic's Q-value for its chosen actions.

Key features:

- Deterministic policy (no entropy term).
- Single critic network (plus its target copy).
- Policy updates come from deterministic policy gradient.

2. Soft Actor-Critic (SAC) Equations

SAC extends the actor-critic framework with maximum entropy RL and twin critics.

Target value (Clipped Double Q-learning):

$$y^{\text{SAC}} = r + \gamma \mathbb{E}_{a' \sim \pi_{\phi}} \left[\min_{i \in \{1,2\}} Q_{\theta_i}(s', a') - \alpha \log \pi_{\phi}(a'|s') \right]$$

Where:

- Two critics $Q_{\theta_1}, Q_{\theta_2}$ — take the minimum to reduce overestimation.
- $\pi_{\phi}(a|s) \rightarrow$ stochastic policy (actor network).
- $\alpha \rightarrow$ temperature parameter controlling entropy weight.
- $-\alpha \log \pi_{\phi}(a'|s') \rightarrow$ entropy bonus (encourages exploration).

Critic loss:

$$L(\theta_i) = \mathbb{E}_{(s,a,r,s') \sim \mathcal{D}} \left[(Q_{\theta_i}(s, a) - y^{\text{SAC}})^2 \right], \quad i = 1, 2$$

Actor loss:

$$J_{\pi}(\phi) = \mathbb{E}_{s \sim \mathcal{D}} [\alpha \log \pi_{\phi}(a|s) - Q_{\theta_1}(s, a)]$$

The actor trades off between high Q-value and high entropy.

Key features:

- Stochastic policy.
- Two critics with minimum selection.
- Explicit entropy regularization in both target and policy objective.

3. Side-by-Side Equation Comparison

Concept	DDPG	SAC
Policy type	Deterministic $\mu_\phi(s)$	Stochastic $\pi_\phi(a s)$
Action in target	$\mu_\phi(s')$	$a' \sim \pi_\phi(a s')$
Critics	One critic (and target)	Two critics (min for stability)
Entropy term	None	$-\alpha \log \pi_\phi(a' s')$
Policy update	Maximize Q-value deterministically	Maximize Q-value + entropy

4. Key Takeaways

- DDPG optimizes a deterministic policy to maximize Q-values, without explicit exploration bonuses.
- SAC optimizes a stochastic policy, balancing exploitation (high Q) and exploration (high entropy).
- SAC's twin-critic and entropy framework make it generally more stable and robust, especially in environments where exploration is important.

Question 10: Policy Gradient Derivation: Fill-in-the-Blank Question

In a one-step MDP (contextual bandit), the agent observes a state $s \sim d(s)$, takes an action $a \sim \pi_\theta(a | s)$, receives a reward $r = R_{s,a}$, and the episode ends.

Let the objective be the expected reward:

$$J(\theta) = \mathbb{E}_{\pi_\theta}[r]$$

Use the score function trick to derive the gradient of $J(\theta)$ with respect to θ .

Question: Complete the missing steps Fill in the missing Line 2 and Line 3 below.

Step	Equation
Line 1	$J(\theta) = \sum_s d(s) \sum_a \pi_\theta(a s) R_{s,a}$
Line 2	<i>You fill this in</i>
Line 3	<i>You fill this in</i>
Line 4	$= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(a s) \cdot r]$

Step	Equation
Line 1	$J(\theta) = \sum_s d(s) \sum_a \pi_\theta(a s) R_{s,a}$
Answer: Line 2	$\nabla_\theta J(\theta) = \sum_s d(s) \sum_a \nabla_\theta \pi_\theta(a s) R_{s,a}$
Line 3	$= \sum_s d(s) \sum_a \pi_\theta(a s) \nabla_\theta \log \pi_\theta(a s) R_{s,a}$
Line 4	$= \mathbb{E}_{\pi_\theta} [\nabla_\theta \log \pi_\theta(a s) \cdot r]$

Refer to section 4.2 for a detailed answer.

Question 11: What is Conservative Q-Learning (CQL) and what problem does it aim to solve in Reinforcement Learning?

Refer to section 17.5 for detailed explanation.

Answer:

- Conservative Q-Learning (CQL) is an *offline reinforcement learning* algorithm designed to learn effective policies from fixed datasets without additional environment interaction.

- The key problem CQL addresses is **overestimation of Q-values** for out-of-distribution (OOD) actions — i.e., actions not well-represented in the offline dataset. Traditional Q-learning methods can assign high values to these unseen actions, leading to poor policy performance when deployed.
- CQL solves this by **penalizing the Q-function for assigning high values to unseen or less-likely actions**. It enforces a form of conservatism by encouraging the learned Q-values to remain low for actions outside the data distribution, thus avoiding optimistic extrapolation.

Mathematically, CQL modifies the Q-learning objective by adding a term like:

$$\text{CQL Objective: } \min_Q \underbrace{\text{Bellman Error}}_{\text{TD Loss}} + \alpha (\mathbb{E}_{a \sim \pi(\cdot|s)}[Q(s, a)] - \mathbb{E}_{a \sim \mathcal{D}}[Q(s, a)])$$

where:

- \mathcal{D} is the dataset distribution,
- π is the learned policy,
- α is a regularization coefficient.

This discourages the policy from assigning high value to actions that are not supported by the dataset.

Question 12: DQN suffers from the problem of overestimation of Q values. The loss function used in DQN is

$$L(w) = \sum_i \left(R_{i+1} + \gamma \max_{a'} \hat{q}_{w^-}(S_{i+1}, a') - \hat{q}_w(S_i, A_i) \right)^2$$

where w^- is the parameters of the target network. Double DQN – an improvement to DQN – mitigates the issue of overestimation in DQN. Double DQN uses the following loss function:

$$L(w) = \sum_i \left(R_{i+1} + \gamma \hat{q}_{w^-} \left(S_{i+1}, \arg \max_{a'} \hat{q}_w(S_{i+1}, a') \right) - \hat{q}_w(S_i, A_i) \right)^2$$

- a. **[2 points]** Explain using the loss function used for training DQN, what causes the overestimation problem in DQN (2 sentences)

Answer:

- $\hat{q}_{w^-}(S_{i+1}, a')$ is a noisy estimate.
- $\max_{a'} \hat{q}_{w^-}(S_{i+1}, a')$ systematically always picks up the max positive noise resulting in overestimation.

- b. **[3 points]** Explain how the loss function used in Double DQN addresses the overestimation problem in DQN (3 sentences)

Answer:

- Double DQN aims to decorrelate the network that selects the max action, and the network that gives the action value corresponding to the selected max action.
- The two networks can have different noisy estimates for action values in a state.
- This reduces the chance of overestimation as for overestimation to happen both networks would need to have high positive direction noise for the same action in a given state.