

# Quiz on Modern Policy Gradient Methods: From First Principles to Soft Actor-Critic

Based on the exposition by Taha Majlesi

July 17, 2025

# Contents

<b>1 Multiple Choice Questions</b>	<b>3</b>
<b>Multiple Choice Questions</b>	<b>3</b>
<b>2 Explainable Questions</b>	<b>22</b>
<b>Explainable Questions</b>	<b>22</b>

# 1 Multiple Choice Questions

**Question 1:** What is the central objective of reinforcement learning as described in the introduction?

- A. To develop an agent that can learn a value function.
- B. To minimize the cumulative reward signal.
- C. To develop an agent that can learn an optimal policy to maximize cumulative reward.
- D. To create a perfect model of the environment's dynamics.

**Answer:**C

**Question 2:** What is the primary challenge in computing the gradient of the objective function  $J(\theta)$  directly?

- A. The reward function is unknown.
- B. The policy parameters  $\theta$  are too complex.
- C. It requires differentiating the environment's dynamics, which are unknown.
- D. The state space is too large.

**Answer:**C

**Question 3:** What mathematical tool is used to make the policy gradient tractable?

- A. The chain rule.
- B. The log-derivative trick.
- C. Taylor expansion.
- D. Fourier analysis.

**Answer:**B

**Question 4:** In the Policy Gradient Theorem, the gradient of the log-probability of a trajectory,  $\nabla_{\theta} \log p_{\theta}(\tau)$ , simplifies to a sum that depends only on:

- A. The environment dynamics.
- B. The reward function.
- C. The policy parameters.
- D. The initial state distribution.

**Answer:**C

**Question 5:** The REINFORCE algorithm is also known as:

- A. Q-Learning.
- B. Deep Deterministic Policy Gradient.
- C. Monte Carlo Policy Gradient.
- D. Temporal Difference Learning.

**Answer:**C

**Question 6:** What is the primary practical drawback of the REINFORCE algorithm?

- A. It is an off-policy algorithm.
- B. It has low bias.
- C. It suffers from high variance in its gradient estimates.
- D. It requires a model of the environment.

**Answer:**C

**Question 7:** How does the "reward-to-go" estimator improve upon using the total return in REINFORCE?

- A. It increases the bias of the gradient estimate.
- B. It respects causality by only considering future rewards.
- C. It uses rewards from the beginning of the episode.
- D. It eliminates variance completely.

**Answer:**B

**Question 8:** What is the purpose of introducing a state-dependent baseline  $b(s_t)$  into the policy gradient formula?

- A. To increase the learning rate.
- B. To reduce the variance of the gradient estimate.
- C. To make the algorithm off-policy.
- D. To introduce bias for faster convergence.

**Answer:**B

**Question 9:** Subtracting a baseline  $b(s_t)$  from the reward-to-go is unbiased as long as the baseline does not depend on:

- A. The state  $s_t$ .
- B. The policy parameters  $\theta$ .
- C. The action  $a_t$ .
- D. The reward  $r_t$ .

**Answer:**C

**Question 10:** What is the ideal choice for the baseline  $b(s_t)$  to achieve optimal variance reduction?

- A. The action-value function  $Q^\pi(s_t, a_t)$ .
- B. The state-value function  $V^\pi(s_t)$ .
- C. A constant value of zero.
- D. The maximum possible reward.

**Answer:**B

**Question 11:** The Advantage Function  $A^\pi(s_t, a_t)$  is defined as:

- A.  $Q^\pi(s_t, a_t) + V^\pi(s_t)$ .
- B.  $V^\pi(s_t) - Q^\pi(s_t, a_t)$ .
- C.  $Q^\pi(s_t, a_t) - V^\pi(s_t)$ .
- D.  $r_t + \gamma V^\pi(s_{t+1})$ .

**Answer:**C

**Question 12:** Generalized Policy Iteration (GPI) involves the interaction of which two processes?

- A. Exploration and Exploitation.
- B. Policy Evaluation and Policy Improvement.
- C. Monte Carlo sampling and Bootstrapping.
- D. On-policy and Off-policy learning.

**Answer:**B

**Question 13:** The Policy Improvement Theorem guarantees that a new policy  $\pi'$  is better than  $\pi$  if for all states  $s$ :

- A.  $V^{\pi'}(s) < V^\pi(s)$ .
- B.  $Q^\pi(s, \pi'(s)) \geq V^\pi(s)$ .

- C.  $Q^\pi(s, \pi'(s)) < V^\pi(s)$ .
- D. The policies are identical.

**Answer:**B

**Question 14:** What does the Performance Difference Lemma express?

- A. The difference in value functions between two states.
- B. The change in performance,  $J(\theta') - J(\theta)$ , in terms of the advantage function.
- C. The error bound of the REINFORCE algorithm.
- D. The variance of the policy gradient estimate.

**Answer:**B

**Question 15:** According to the Performance Difference Lemma, the performance improvement is an expectation over trajectories sampled from:

- A. The old policy  $\pi_\theta$ .
- B. A uniform random policy.
- C. The new policy  $\pi_{\theta'}$ .
- D. The optimal policy  $\pi^*$ .

**Answer:**C

**Question 16:** Policy gradient methods can be interpreted as an "infinitesimal" version of Policy Iteration because:

- A. They take large, discrete steps.
- B. They only work in discrete action spaces.
- C. They take small steps in the direction of the gradient.
- D. They do not use a value function.

**Answer:**C

**Question 17:** What is the "state distribution mismatch" problem in on-policy learning?

- A. The initial state distribution is unknown.
- B. The state distribution changes too slowly.
- C. Data from the old policy is used to evaluate the new policy.
- D. The state distribution is not a normal distribution.

**Answer:**C

**Question 18:** The on-policy approximation for the performance improvement is only valid when:

- A. The learning rate is very high.
- B. The policy updates are kept small.
- C. The environment is deterministic.
- D. The reward function is always positive.

**Answer:**B

**Question 19:** What is the statistically sound method for correcting distribution mismatch?

- A. Adding a baseline.
- B. Using a smaller neural network.
- C. Importance Sampling (IS).
- D. Increasing the batch size.

**Answer:**C

**Question 20:** Why is a full off-policy correction using importance sampling often practically unusable?

- A. It introduces significant bias.
- B. It suffers from extremely high variance.
- C. It is computationally too simple.
- D. It only works for deterministic policies.

**Answer:**B

**Question 21:** The theoretical bound on the state distribution mismatch shows that the divergence between state distributions grows in proportion to:

- A. The discount factor  $\gamma$ .
- B. The number of states.
- C. The "closeness" of the policies,  $\epsilon$ , and time,  $t$ .
- D. The magnitude of the rewards.

**Answer:**C

**Question 22:** Maximizing the approximate objective  $L(\theta')$  is equivalent to maximizing what?

- A. The variance of the gradient.
- B. An upper bound on the true performance.
- C. A lower bound on the true performance improvement.
- D. The KL-divergence between policies.

**Answer:**C

**Question 23:** What is the core idea behind Trust Region methods like TRPO and PPO?

- A. To maximize the size of policy updates.
- B. To explicitly manage the trade-off between learning speed and stability.
- C. To eliminate the need for a value function.
- D. To use a deterministic policy.

**Answer:**B

**Question 24:** How does TRPO ensure that the new policy stays "close" to the old one?

- A. By using a very small learning rate.
- B. By clipping the objective function.
- C. By solving a constrained optimization problem with a KL-divergence constraint.
- D. By using a replay buffer.

**Answer:**C

**Question 25:** What is a major drawback of TRPO?

- A. It is unstable.
- B. It is computationally expensive and complex to implement.
- C. It has very low sample efficiency.
- D. It cannot handle continuous action spaces.

**Answer:**B

**Question 26:** How does PPO approximate the behavior of TRPO?

- A. By using a hard KL-divergence constraint.



- B. By using a clipped surrogate objective function.
- C. By using second-order optimization.
- D. By using a deterministic policy.

**Answer:**B

**Question 27:** The clipped surrogate objective in PPO effectively creates a:

- A. Hard constraint.
- B. Soft "trust region".
- C. Deterministic policy update.
- D. Larger state space.

**Answer:**B

**Question 28:** What is the new paradigm introduced by Soft Actor-Critic (SAC)?

- A. Minimum entropy reinforcement learning.
- B. On-policy actor-critic learning.
- C. Model-based reinforcement learning.
- D. Maximum entropy reinforcement learning.

**Answer:**D

**Question 29:** What is the main benefit of adding an entropy term to the RL objective?

- A. It simplifies the algorithm.
- B. It encourages exploration and improves robustness.
- C. It guarantees convergence to the optimal policy in all cases.
- D. It reduces the computational cost.

**Answer:**B

**Question 30:** The SAC algorithm is derived from a modified framework called:

- A. Hard Policy Iteration.
- B. Generalized Policy Iteration.
- C. Soft Policy Iteration.
- D. Value Iteration.

**Answer:**C

**Question 31:** The soft Q-function in SAC is defined as the expected return plus:

- A. The expected future discounted rewards only.
- B. The expected future discounted entropy.
- C. The variance of the policy.
- D. The KL-divergence from a prior policy.

**Answer:**B

**Question 32:** In Soft Policy Improvement, the new policy is updated to be close to what distribution?

- A. A uniform distribution.
- B. A greedy distribution (argmax).
- C. A Gaussian distribution centered at zero.
- D. The Boltzmann distribution defined by the soft Q-function.

**Answer:**D

**Question 33:** What does Lemma 2 (Soft Policy Improvement) in the SAC theory guarantee?

- A. That the policy entropy will always decrease.
- B. Monotonic improvement in the soft Q-value.
- C. That the algorithm will converge in a single step.
- D. That the policy will become deterministic.

**Answer:**B

**Question 34:** SAC is what type of algorithm?

- A. On-policy, value-based.
- B. Off-policy, actor-critic.
- C. On-policy, actor-critic.
- D. Off-policy, value-based.

**Answer:**B

**Question 35:** Why does SAC use two Q-Critic networks?

- A. To speed up computation.
- B. To handle two different reward signals.
- C. To combat the overestimation of Q-values (Clipped Double-Q Learning).
- D. To have one for exploration and one for exploitation.

**Answer:**C

**Question 36:** What is the purpose of the target networks in SAC?

- A. To provide a stable learning target for the critic updates.
- B. To learn the policy directly.
- C. To estimate the reward function.
- D. To increase the variance of the updates.

**Answer:**A

**Question 37:** How are the target networks in SAC typically updated?

- A. By direct copy of the main network weights every step.
- B. By Polyak (soft) averaging.
- C. By retraining from scratch periodically.
- D. They are not updated.

**Answer:**B

**Question 38:** What technique does SAC use to allow gradients to be backpropagated through a stochastic policy?

- A. The log-derivative trick.
- B. The reparameterization trick.
- C. Monte Carlo sampling.
- D. Importance sampling.

**Answer:**B

**Question 39:** What is the purpose of learning the temperature parameter  $\alpha$  in SAC?

- A. To keep the entropy constant at zero.
- B. To automate the trade-off between exploration and exploitation.
- C. To increase the learning rate.

D. To simplify the Q-function loss.

**Answer:**B

**Question 40:** Compared to on-policy algorithms like PPO, SAC is generally considered to have higher:

- A. Stability.
- B. Simplicity.
- C. Sample efficiency.
- D. Variance.

**Answer:**C

**Question 41:** The objective function  $J(\theta)$  in policy gradient methods typically represents the expected:

- A. Immediate reward.
- B. Total discounted reward.
- C. Policy entropy.
- D. State visitation frequency.

**Answer:**B

**Question 42:** The final form of the policy gradient is an expectation over trajectories. How is this expectation estimated in practice?

- A. By solving a system of linear equations.
- B. By averaging over a batch of sampled trajectories.
- C. By using a closed-form analytical solution.
- D. By differentiating a model of the environment.

**Answer:**B

**Question 43:** In the context of variance reduction, what does the advantage function  $A^\pi(s, a)$  measure?

- A. The absolute goodness of an action.
- B. How much better an action is compared to the policy's average behavior in that state.
- C. The probability of reaching a terminal state.
- D. The total reward from the start of the episode.

**Answer:**B

**Question 44:** Actor-critic methods use a "critic" to estimate:

- A. The policy parameters.
- B. The environment's dynamics.
- C. The value function (or advantage function).
- D. The reward function.

**Answer:**C

**Question 45:** The "policy improvement" step in GPI aims to find a new policy that is:

- A. More random.
- B. "Greedy" with respect to the current value function.
- C. Easier to compute.
- D. A weighted average of all previous policies.

**Answer:**B

**Question 46:** The term inside the summation of the Performance Difference Lemma,  $r(s_t, a_t) + \gamma V^{\pi_\theta}(s_{t+1}) - V^{\pi_\theta}(s_t)$ , is also known as the:

- A. Monte Carlo return.
- B. State visitation frequency.
- C. Temporal-Difference (TD) error.
- D. Policy entropy.

**Answer:**C

**Question 47:** A large policy update can lead to "policy collapse" because:

- A. The learning rate becomes too small.
- B. The value function estimate becomes perfect.
- C. The on-policy data approximation becomes invalid, leading to a performance drop.
- D. The neural network runs out of memory.

**Answer:**C

**Question 48:** The importance ratio  $r_t(\theta)$  in PPO is defined as:

- A.  $\pi_\theta(a_t|s_t) - \pi_{\theta_{\text{old}}}(a_t|s_t)$ .

- B.  $\frac{\pi_{\theta_{\text{old}}}(a_t|s_t)}{\pi_{\theta}(a_t|s_t)}$ .
- C.  $\frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}$ .
- D.  $\log \pi_{\theta}(a_t|s_t)$ .

**Answer:**C

**Question 49:** PPO is generally preferred over TRPO in practice because it is:

- A. A second-order method.
- B. More theoretically sound.
- C. Simpler to implement and more sample-efficient.
- D. Guaranteed to find the global optimum.

**Answer:**C

**Question 50:** A maximum entropy policy is beneficial because it is less "brittle," meaning it:

- A. Always chooses the single best action.
- B. Does not commit fully to one course of action and can adapt.
- C. Converges faster than any other method.
- D. Uses a linear function approximator.

**Answer:**B

**Question 51:** The proof of convergence for Soft Policy Iteration relies on the fact that the sequence of soft Q-functions is:

- A. Monotonically decreasing and bounded.
- B. Oscillating around the optimal value.
- C. Monotonically increasing and bounded.
- D. Unbounded.

**Answer:**C

**Question 52:** In the SAC Q-function loss, the target value  $y$  includes a term  $-\alpha \log \pi_{\phi}(a_{t+1}|s_{t+1})$ . This term corresponds to the:

- A. Advantage function.
- B. Reward signal.
- C. Entropy bonus.

D. Importance sampling ratio.

**Answer:**C

**Question 53:** The policy loss in SAC,  $J_\pi(\phi)$ , aims to find actions that maximize the:

- A. Policy entropy only.
- B. Soft Q-function value.
- C. Distance from the previous policy.
- D. Reward at the current timestep.

**Answer:**B

**Question 54:** The evolution from REINFORCE to SAC shows a trend towards addressing which key challenges?

- A. Model accuracy, state representation, and reward design.
- B. Variance, stability, and sample efficiency.
- C. Hardware limitations, memory usage, and computational speed.
- D. Discrete actions, small state spaces, and linear models.

**Answer:**B

**Question 55:** The term  $\nabla_\theta \log \pi_\theta(a_t|s_t)$  is often called the:

- A. Value function.
- B. Score function.
- C. Reward function.
- D. Transition function.

**Answer:**B

**Question 56:** If all rewards in an environment are shifted by a constant  $C$ , how does this affect the policy gradient if a proper baseline is used?

- A. The gradient becomes zero.
- B. The learning process becomes unstable.
- C. The gradient is largely unaffected because the advantage is centered.
- D. The gradient is also shifted by  $C$ .

**Answer:**C

**Question 57:** The "actor" in an actor-critic method corresponds to the:

- A. Value function.
- B. Policy.
- C. Environment model.
- D. Reward model.

**Answer:**B

**Question 58:** The "generalized" aspect of GPI refers to the fact that evaluation and improvement steps can be:

- A. Skipped entirely.
- B. Performed only once.
- C. Interleaved in many ways.
- D. Done in parallel on different machines.

**Answer:**C

**Question 59:** The use of an experience replay buffer is a key feature of which type of algorithms?

- A. On-policy algorithms.
- B. Off-policy algorithms.
- C. Model-based algorithms.
- D. Gradient-free algorithms.

**Answer:**B

**Question 60:** In the PPO clipped objective, what happens if the advantage  $A_t$  is negative?

- A. The objective is maximized when the ratio  $r_t(\theta)$  is large.
- B. The clipping encourages making the ratio  $r_t(\theta)$  smaller (i.e., making the action less likely).
- C. The update is skipped for that timestep.
- D. The advantage is ignored.

**Answer:**B

**Question 61:** SAC's ability to learn multi-modal behaviors means it can:

- A. Only solve problems with one optimal solution.



- B. Learn to assign probability to multiple, equally good strategies.
- C. Only work in environments with visual inputs.
- D. Switch between on-policy and off-policy learning.

**Answer:**B

**Question 62:** The soft Bellman backup operator  $\mathcal{T}^\pi$  is proven to be a:

- A. Expansion mapping.
- B. Contraction mapping.
- C. Identity mapping.
- D. Zero mapping.

**Answer:**B

**Question 63:** The target entropy  $\bar{\mathcal{H}}$  in the automatic temperature adjustment of SAC is a:

- A. Learned parameter.
- B. Value that changes every timestep.
- C. Hyperparameter set by the user.
- D. Function of the current state.

**Answer:**C

**Question 64:** Which algorithm from the list is NOT an on-policy algorithm?

- A. REINFORCE.
- B. TRPO.
- C. PPO.
- D. SAC.

**Answer:**D

**Question 65:** The derivation of the Policy Gradient Theorem crucially relies on the fact that the gradient of the environment dynamics with respect to policy parameters is:

- A. One.
- B. Infinite.
- C. Zero.
- D. Unknown but non-zero.

**Answer:**C

**Question 66:** A policy gradient update  $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$  is an application of:

- A. Gradient descent.
- B. Gradient ascent.
- C. Newton's method.
- D. Conjugate gradient.

**Answer:**B

**Question 67:** The TD error can be seen as a single-sample estimate of the:

- A. Total return.
- B. Advantage function.
- C. Policy entropy.
- D. State distribution.

**Answer:**B

**Question 68:** TRPO is considered a second-order optimization method because it uses information related to the:

- A. Hessian (or Fisher Information Matrix).
- B. Gradient only.
- C. Learning rate schedule.
- D. Batch size.

**Answer:**A

**Question 69:** The clipping parameter  $\epsilon$  in PPO controls the size of the:

- A. Neural network.
- B. Replay buffer.
- C. Soft trust region.
- D. Discount factor.

**Answer:**C

**Question 70:** The "temperature" parameter  $\alpha$  in SAC controls the trade-off between reward and:

- A. Bias.
- B. Variance.
- C. Entropy.
- D. Stability.

**Answer:**C

**Question 71:** The reparameterization trick works by reformulating the sampling of an action  $a \sim \pi_\phi(\cdot|s)$  as:

- A. A discrete choice from a set.
- B. A deterministic function of state and a random noise vector.
- C. A lookup in a large table.
- D. A query to a model of the environment.

**Answer:**B

**Question 72:** Which algorithm introduced the idea of a clipped surrogate objective?

- A. REINFORCE.
- B. TRPO.
- C. PPO.
- D. SAC.

**Answer:**C

**Question 73:** The intellectual lineage of policy gradient methods shows a progression from solving the problem of high variance to solving the problem of:

- A. Low bias.
- B. Learning instability due to large updates.
- C. Inaccurate reward functions.
- D. Discrete state spaces.

**Answer:**B

**Question 74:** The term  $p(s_{t+1}|s_t, a_t)$  represents the:

- A. Policy.
- B. Reward function.
- C. Environment's transition dynamics.

D. Value function.

**Answer:**C

**Question 75:** In the context of GPI, "policy evaluation" aims to answer the question:

- A. "What is the best action to take?"
- B. "How can I improve my policy?"
- C. "How good is the current policy?"
- D. "What will the next state be?"

**Answer:**C

**Question 76:** The proof that subtracting a baseline is unbiased relies on the fact that  $\nabla_{\theta} \int \pi_{\theta}(a|s) da = \nabla_{\theta}(1)$ , which equals:

- A. 1.
- B. 0.
- C.  $\infty$ .
- D.  $\pi_{\theta}(a|s)$ .

**Answer:**B

**Question 77:** The Clipped Double-Q Learning trick used in SAC is borrowed from which algorithm?

- A. DDPG.
- B. A3C.
- C. PPO.
- D. TD3.

**Answer:**D

**Question 78:** A key benefit of SAC being an off-policy algorithm is its:

- A. High stability.
- B. Low variance.
- C. High sample efficiency from reusing old data.
- D. Simplicity of implementation.

**Answer:**C

**Question 79:** The update for the policy in PPO is performed using:

- A. Second-order optimization.
- B. First-order optimization (gradient ascent).
- C. A closed-form solution.
- D. A genetic algorithm.

**Answer:**B

**Question 80:** The final conclusion of the report is that the evolution from REINFORCE to SAC is a narrative of overcoming challenges in:

- A. Hardware, software, and networking.
- B. Supervised, unsupervised, and semi-supervised learning.
- C. Variance, stability, and sample efficiency.
- D. Feature engineering, data collection, and model selection.

**Answer:**C

## 2 Explainable Questions

**Question 1:** Explain the log-derivative trick and why it is so crucial for deriving the Policy Gradient Theorem.

**Explanation:** The log-derivative trick is a mathematical identity that states  $\nabla f(x) = f(x)\nabla \log f(x)$ . It is crucial for the Policy Gradient Theorem because it allows us to solve a seemingly intractable problem.

The initial form of the policy gradient is  $\nabla_{\theta} J(\theta) = \int [\nabla_{\theta} p_{\theta}(\tau)] R(\tau) d\tau$ . The term  $\nabla_{\theta} p_{\theta}(\tau)$  is problematic because the trajectory distribution  $p_{\theta}(\tau)$  depends on the environment's dynamics,  $p(s_{t+1}|s_t, a_t)$ , which are unknown. Therefore, we cannot differentiate it directly.

By applying the log-derivative trick, we transform the problematic term:

$$\nabla_{\theta} p_{\theta}(\tau) = p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau)$$

Substituting this back into the gradient expression gives:

$$\nabla_{\theta} J(\theta) = \int p_{\theta}(\tau) [\nabla_{\theta} \log p_{\theta}(\tau)] R(\tau) d\tau$$

This expression is now an expectation under the trajectory distribution  $p_{\theta}(\tau)$ :

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}} [(\nabla_{\theta} \log p_{\theta}(\tau)) R(\tau)]$$

The final, crucial step is that when we expand  $\nabla_{\theta} \log p_{\theta}(\tau)$ , the terms related to the environment dynamics have a gradient of zero with respect to  $\theta$  and drop out. This leaves a gradient that depends only on the policy,  $\pi_{\theta}$ , which the agent controls. The trick thus transforms an intractable derivative into a tractable expectation that can be estimated by sampling.

**Question 2:** Describe the problem of high variance in the REINFORCE algorithm and explain intuitively and mathematically how introducing a state-dependent baseline helps mitigate it.

**Explanation:** The problem of high variance in REINFORCE stems from its gradient estimator, which scales the score function  $\nabla_{\theta} \log \pi_{\theta}(a_t|s_t)$  by the Monte Carlo return (reward-to-go)  $\hat{Q}_t^{\pi}$ . Because both the policy and the environment can be stochastic, the return  $\hat{Q}_t^{\pi}$  for the same state-action pair can vary significantly across different episodes. This makes the gradient estimate very noisy, leading to unstable updates and slow learning.

**Intuitive Explanation:** Imagine an environment where all rewards are large and positive (e.g., between 100 and 110). An action leading to a return of 101 and an action leading to 109 will both be reinforced (pushed to be more likely), even though the latter is clearly better. The absolute magnitude of the return doesn't provide a good learning signal. What really matters is whether an action was *better or worse than average* for a given state. A baseline provides this "average" reference point. By subtracting the baseline, we center the learning signal. Actions better than average get a positive push, and actions worse than average get a negative push, which is a much more informative and stable signal.

**Mathematical Explanation:** We introduce a baseline  $b(s_t)$  that depends only on the state, creating the new gradient estimator:

$$\nabla_{\theta} J(\theta) = \left[ \sum_t \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) (\hat{Q}_t^{\pi} - b(s_t)) \right]$$

This estimator remains unbiased because the expected value of the added term is zero:  $[\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) b(s_t)] = 0$ . The variance of this new estimator is minimized when the baseline  $b(s_t)$  is highly correlated with  $\hat{Q}_t^{\pi}$ . The optimal choice for the baseline is the true state-value function,  $b(s_t) = V^{\pi}(s_t) = [\hat{Q}_t^{\pi} | s_t]$ . This choice transforms the learning signal into the **Advantage Function**,  $A^{\pi}(s_t, a_t) = Q^{\pi}(s_t, a_t) - V^{\pi}(s_t)$ , which has lower variance because it subtracts the expected return from the sampled return.

**Question 3:** Explain the core components of Generalized Policy Iteration (GPI) and how the policy gradient update process can be interpreted as an approximate form of GPI.

**Explanation:** Generalized Policy Iteration (GPI) is a general framework that describes the interaction of two processes:

1. **Policy Evaluation:** Given a policy  $\pi$ , compute its value function,  $V^{\pi}$  or  $Q^{\pi}$ . This step answers the question, "How good is my current strategy?"
2. **Policy Improvement:** Given a value function, update the policy to be "greedy" with respect to it, creating a new, better policy  $\pi'$ . This step answers the question, "How can I act better based on my evaluation?"

The interaction of these two processes is guaranteed to converge to the optimal policy and value function.

The policy gradient process can be interpreted as an approximate form of GPI:

1. **Approximate Policy Evaluation:** Instead of computing the exact value function, policy gradient methods estimate the **advantage function**,  $A^{\pi}(s, a)$ . In REINFORCE, this is done with Monte Carlo rollouts. In actor-critic methods, a "critic" network is trained to approximate the value function, which is then used to compute an advantage estimate. This is an "approximate" evaluation because it uses sampled data and function approximators.
2. **Infinitesimal Policy Improvement:** Instead of making a discrete, greedy update, policy gradient methods take a small step in the direction that improves the policy. The update rule,  $\theta \leftarrow \theta + \alpha [\nabla_{\theta} \log \pi_{\theta}(a | s) A^{\pi}(s, a)]$ , pushes the policy to increase the probability of actions with positive advantages and decrease the probability of those with negative advantages. This is analogous to the policy improvement step, but it is "infinitesimal" because it's a small gradient step, not a large, definitive change.

This perspective is formalized by the Performance Difference Lemma, which shows that the expected advantage is precisely the quantity that determines the improvement in policy performance, linking the gradient update directly to the core principle of policy improvement in GPI.

**Question 4:** What is the state distribution mismatch problem, and why does it necessitate trust region methods like TRPO and PPO?

**Explanation:** The state distribution mismatch problem is a subtle but critical flaw in the practical application of on-policy policy gradient methods.

The theory (specifically, the Performance Difference Lemma) states that the improvement in a policy,  $J(\theta') - J(\theta)$ , is the expected advantage of the *new* policy's actions, evaluated under the state distribution induced by the *new* policy,  $\pi_{\theta'}$ . However, to compute this update, we only have data (trajectories) collected using the *old* policy,  $\pi_{\theta}$ .

On-policy algorithms make the practical approximation of using the old data, effectively assuming that the state distribution of the new policy is the same as the old one ( $p_{\theta'}(s) \approx p_{\theta}(s)$ ). This approximation is only reasonable if the new policy is very similar to the old one. If we take a large update step,  $\pi_{\theta'}$  might be very different from  $\pi_{\theta}$ , causing it to visit a completely different set of states. In this case, the data from the old policy is no longer a valid indicator of the new policy's performance, and the update can be catastrophic, leading to a sharp drop in performance (policy collapse).

This creates a fundamental tension: we want large updates for fast learning, but large updates can be unstable. **Trust region methods** are designed to solve this exact problem.

- **TRPO** formalizes this by adding a hard constraint to the optimization: it maximizes the performance objective while explicitly constraining the KL-divergence between the old and new policies to be within a small radius  $\delta$  (the "trust region"). This guarantees that the new policy is close to the old one, ensuring the approximation is valid and leading to monotonic improvement.
- **PPO** achieves a similar goal in a simpler, first-order way. Its clipped surrogate objective function removes the incentive for making the policy ratio  $r_t(\theta)$  too large or too small. This effectively creates a soft trust region that discourages overly large policy updates, providing stability without the computational complexity of TRPO.

**Question 5:** Explain the maximum entropy objective in Soft Actor-Critic (SAC) and discuss its three main benefits.

**Explanation:** The standard reinforcement learning objective is to maximize the cumulative reward. Soft Actor-Critic (SAC) modifies this by adding a policy entropy term, creating the **maximum entropy objective**:

$$J(\pi) = \sum_{t=0}^T \mathbb{E}_{(s_t, a_t) \sim \rho_{\pi}} [r(s_t, a_t) + \alpha \mathcal{H}(\pi(\cdot | s_t))]$$

Here,  $\mathcal{H}(\pi(\cdot | s_t))$  is the entropy (randomness) of the policy at state  $s_t$ , and  $\alpha$  is a temperature parameter that balances the reward and entropy terms. The agent is thus incentivized not only to succeed at the task (maximize reward) but also to act as randomly as possible while doing so.

This objective has three main benefits:

1. **Improved Exploration:** The entropy bonus explicitly encourages the agent to try different actions. By rewarding randomness, it prevents the policy from prematurely



converging to a single, sharp, and possibly suboptimal strategy. This allows the agent to explore more of the state-action space and discover more diverse and potentially better solutions.

2. **Improved Robustness and Stability:** Policies learned with maximum entropy are less "brittle." Since they don't commit fully to a single action (unless it's vastly superior), they can adapt more easily if the environment changes or if some actions become unavailable. For example, if two actions have similar high Q-values, the SAC policy will assign them roughly equal probability. This also helps to smooth the learning landscape, leading to more stable convergence during training.
3. **Learning Multi-modal Behaviors:** In tasks where multiple, distinct strategies are equally optimal, a standard RL agent might arbitrarily converge to just one. A maximum entropy agent, however, can learn to represent all of them by assigning probability mass to each optimal mode of behavior, resulting in a more versatile policy.

**Question 6:** Describe the practical architecture of the SAC algorithm, detailing the roles of the twin Q-critics, target networks, and the reparameterization trick.

**Explanation:** The practical SAC algorithm is an off-policy actor-critic method with several key architectural components designed for stability and efficiency.

#### Architecture Overview:

- **Actor ( $\pi_\phi$ ):** A stochastic policy network that outputs a probability distribution over actions for a given state.
- **Twin Q-Critics ( $Q_{\theta_1}, Q_{\theta_2}$ ):** Two separate Q-function networks that estimate the soft Q-value for a state-action pair.
- **Twin Target Q-Networks ( $Q_{\theta_{1,\text{targ}}}, Q_{\theta_{2,\text{targ}}}$ ):** Time-delayed copies of the critic networks, used to provide stable learning targets.
- **Experience Replay Buffer ( $\mathcal{D}$ ):** Stores past transitions  $(s, a, r, s')$  for off-policy learning.

#### Roles of Key Components:

1. **Twin Q-Critics (Clipped Double-Q Learning):** Standard Q-learning methods are prone to overestimating Q-values, which can lead to suboptimal policies. SAC uses two independent Q-critics and, when calculating the target value for the Bellman update, it uses the *minimum* of the two target Q-network predictions:

$$y = r_t + \gamma(\min_{j=1,2} Q_{\theta_{j,\text{targ}}}(s_{t+1}, a_{t+1}) - \dots)$$

This technique, borrowed from TD3, helps to combat this overestimation bias, leading to more stable and reliable learning of the Q-function.

2. **Target Networks:** The Q-function loss is calculated as the difference between the current Q-estimate and a target value  $y$ . If this target value were calculated using the same Q-networks that are being actively trained, the learning target would be

constantly shifting, which can lead to oscillations and instability. Target networks solve this by providing a stable target. Their weights are updated slowly, typically via Polyak averaging ( $\theta_{\text{targ}} \leftarrow \rho\theta_{\text{targ}} + (1 - \rho)\theta$ ), creating a delayed but more stable version of the critic.

3. **Reparameterization Trick:** The policy needs to be updated by backpropagating a gradient from the Q-function. However, you cannot backpropagate through a random sampling process. The reparameterization trick makes the policy gradient tractable. Instead of sampling an action  $a$  directly from the policy distribution  $\pi(\cdot|s)$ , the action is formulated as a deterministic function of the state and a random noise vector  $\epsilon$  sampled from a fixed distribution (e.g., a standard Gaussian). For a Gaussian policy, this looks like:

$$a_t = \mu_\phi(s_t) + \sigma_\phi(s_t) \odot \epsilon_t, \quad \epsilon_t \sim \mathcal{N}(0, I)$$

Now, the stochasticity is external to the network, and the gradient can flow through the deterministic function  $f_\phi(\epsilon_t; s_t)$  (i.e., through  $\mu_\phi$  and  $\sigma_\phi$ ) back to the policy parameters  $\phi$ . This provides a low-variance way to train the stochastic actor.

**Question 7:** Compare and contrast TRPO and PPO. Why is PPO often preferred in practice?

**Explanation:** TRPO (Trust Region Policy Optimization) and PPO (Proximal Policy Optimization) are both on-policy algorithms designed to solve the stability problem in policy gradient methods by constraining the size of policy updates.

**Similarities:**

- **Goal:** Both aim to make the largest possible policy improvement step without causing a performance collapse, effectively keeping the new policy within a "trust region" of the old one.
- **On-Policy:** Both are on-policy methods, meaning they learn from data collected with the most recent version of the policy.
- **Foundation:** Both are based on the same theoretical motivation: maximizing the performance improvement objective, which is approximated using data from the old policy.

**Differences:**

Feature	TRPO	PPO
<b>Constraint</b>	Hard constraint	Soft constraint
<b>Mechanism</b>	Solves a constrained optimization problem with a strict KL-divergence limit ( $\leq \delta$ ).	Uses a clipped surrogate objective function that penalizes large policy ratios.
<b>Optimization</b>	Second-order method	First-order method
<b>Complexity</b>	High. Requires complex calculations like the Fisher-vector product and conjugate gradient descent.	Low. Implemented with standard stochastic gradient ascent.

**Why PPO is Preferred in Practice:** PPO is often preferred over TRPO for several key reasons, all stemming from its simplicity:

1. **Simplicity of Implementation:** TRPO is notoriously difficult to implement correctly due to its second-order optimization machinery. PPO's clipped objective is much simpler and can be implemented with only a few lines of code in standard deep learning frameworks.
2. **Computational Efficiency:** As a first-order method, PPO is computationally cheaper than TRPO. It avoids the expensive per-step calculations of TRPO, allowing for faster iterations.
3. **Performance:** Despite its simplicity, PPO has been shown empirically to achieve comparable or even better performance than TRPO across a wide range of benchmark tasks. It strikes an excellent balance between stability, performance, and ease of use.

In essence, PPO provides most of the stability benefits of TRPO without the implementation and computational overhead, making it a more practical and widely used default choice for on-policy reinforcement learning.

**Question 8:** Derive the unbiased nature of the baseline in the policy gradient estimator.

**Explanation:** The goal is to show that subtracting a state-dependent baseline  $b(s_t)$  from the policy gradient estimator does not introduce bias. This means we need to prove that the expectation of the term we add is zero. The term added to the gradient objective is  $-\nabla_{\theta} \log \pi_{\theta}(a_t|s_t)b(s_t)$ . We want to show:

$$\mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_{t=0}^T \nabla_{\theta} \log \pi_{\theta}(a_t|s_t)b(s_t) \right] = 0$$

We can analyze the expectation at a single timestep  $t$ , as the expectation of the sum is the sum of expectations. The expectation is over states  $s_t$  and actions  $a_t$  determined by the policy  $\pi_{\theta}$ .

$$\mathbb{E}_{(s_t, a_t) \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a_t|s_t)b(s_t)]$$

We can break this down using the law of iterated expectations:  $\mathbb{E}[X] = \mathbb{E}_Y[\mathbb{E}[X|Y]]$ . Here, we first take the expectation over actions  $a_t$  given the state  $s_t$ , and then over states  $s_t$ .

$$= \mathbb{E}_{s_t \sim \pi_{\theta}} [\mathbb{E}_{a_t \sim \pi_{\theta}(\cdot|s_t)} [\nabla_{\theta} \log \pi_{\theta}(a_t|s_t)b(s_t)|s_t]]$$

Since  $b(s_t)$  is constant with respect to the action  $a_t$  at state  $s_t$ , we can pull it out of the inner expectation:

$$= \mathbb{E}_{s_t \sim \pi_{\theta}} [b(s_t) \mathbb{E}_{a_t \sim \pi_{\theta}(\cdot|s_t)} [\nabla_{\theta} \log \pi_{\theta}(a_t|s_t)]]$$

Now, let's focus on the inner expectation. We can write it as an integral (or sum for discrete actions):

$$\mathbb{E}_{a_t \sim \pi_{\theta}(\cdot|s_t)} [\nabla_{\theta} \log \pi_{\theta}(a_t|s_t)] = \int \pi_{\theta}(a_t|s_t) \nabla_{\theta} \log \pi_{\theta}(a_t|s_t) da_t$$

Using the log-derivative trick ( $\nabla f(x) = f(x) \nabla \log f(x)$ ) in reverse, we get:

$$= \int \nabla_{\theta} \pi_{\theta}(a_t|s_t) da_t$$

We can swap the gradient and the integral (assuming mild regularity conditions):

$$= \nabla_{\theta} \int \pi_{\theta}(a_t|s_t) da_t$$

The integral of a probability distribution over its entire domain is, by definition, equal to 1.

$$= \nabla_{\theta}(1)$$

The gradient of a constant is zero.

$$= 0$$

Substituting this result back, the entire expression becomes:

$$_{s_t \sim \pi_{\theta}} [b(s_t) \cdot 0] = 0$$

Since the expected value of the added term is zero for every timestep  $t$ , the expectation of the sum is also zero. Therefore, subtracting a state-dependent baseline does not change the expected gradient, and the estimator remains unbiased.

**Question 9:** What is the Soft Bellman Equation, and how does it differ from the standard Bellman Equation?

**Explanation:** The Bellman equations are the foundation of value-based reinforcement learning, relating the value of a state to the values of successor states.

The **standard Bellman expectation equation** for the action-value function  $Q^{\pi}(s, a)$  is:

$$Q^{\pi}(s_t, a_t) = [r_{t+1} + \gamma Q^{\pi}(s_{t+1}, a_{t+1})]$$

This can be written as:

$$Q^{\pi}(s_t, a_t) = r(s_t, a_t) + \gamma_{s_{t+1} \sim p, a_{t+1} \sim \pi} [Q^{\pi}(s_{t+1}, a_{t+1})]$$

It states that the value of taking action  $a_t$  in state  $s_t$  is the immediate reward plus the discounted expected value of the next state-action pair, following policy  $\pi$ .

The **Soft Bellman expectation equation**, which is the foundation of Soft Actor-Critic, incorporates the maximum entropy objective. It modifies the standard equation by including the policy entropy at the next step. The soft action-value function  $Q^{\pi}(s, a)$  is defined by:

$$Q^{\pi}(s_t, a_t) = r(s_t, a_t) + \gamma_{s_{t+1} \sim p, a_{t+1} \sim \pi} [Q^{\pi}(s_{t+1}, a_{t+1}) - \alpha \log \pi(a_{t+1}|s_{t+1})]$$

**Key Differences:**

1. **Entropy Term:** The most obvious difference is the presence of the entropy term,  $-\alpha \log \pi(a_{t+1}|s_{t+1})$ . The standard Bellman equation only considers future rewards. The soft version considers both future rewards and future entropy.
2. **Definition of Value:** In the standard framework, the value function represents the expected cumulative reward. In the soft framework, the value function represents the expected cumulative *reward plus entropy*.

3. **Optimal Policy:** The optimal policy derived from the standard Bellman optimality equation is greedy and often deterministic. The optimal policy derived from the soft Bellman optimality equation is stochastic (unless the Q-values are sharply peaked) and balances maximizing rewards with maintaining high entropy.

In essence, the Soft Bellman Equation redefines "value" to include not just the utility of rewards but also the utility of acting randomly and unpredictably.

**Question 10:** Explain the concept of "policy collapse" and how it relates to the state distribution mismatch problem.

**Explanation:** **Policy collapse** is a catastrophic failure mode in on-policy reinforcement learning where a single policy update leads to a sudden and dramatic drop in the agent's performance, from which it may not be able to recover.

The phenomenon is a direct consequence of the **state distribution mismatch problem**. Here's the chain of events:

1. **Data Collection:** The agent uses its current policy,  $\pi_\theta$ , to collect a batch of experience (trajectories).
2. **Gradient Estimation:** It uses this data to estimate the policy gradient. The core assumption here is that this data is a good representation of how a slightly improved policy,  $\pi_{\theta'}$ , will perform.
3. **Large Update Step:** Due to stochasticity in the gradient estimate or a large learning rate, the agent takes a very large step in parameter space, resulting in a new policy  $\pi_{\theta'}$  that is significantly different from the old policy  $\pi_\theta$ .
4. **Approximation Failure:** Because  $\pi_{\theta'}$  is now very different from  $\pi_\theta$ , the state distribution it induces,  $p_{\theta'}(s)$ , is also very different from the old distribution,  $p_\theta(s)$ . The fundamental approximation that the old data is representative of the new policy's performance is now invalid. The update was calculated based on the promise of high rewards in states visited by  $\pi_\theta$ , but the new policy  $\pi_{\theta'}$  may now primarily visit a completely different set of states where its performance is actually terrible.
5. **Performance Collapse:** When the new policy  $\pi_{\theta'}$  is executed in the environment, its actual performance is found to be much worse than what the update predicted. The agent has essentially updated itself into a poor region of the policy space based on misleading information.

In short, policy collapse happens when an overly aggressive update step breaks the on-policy learning assumption, leading the agent to trust a gradient estimate that is no longer relevant for the resulting new policy. This is precisely the problem that trust region methods like TRPO and PPO are designed to prevent by explicitly limiting the magnitude of the policy update.

**Question 11:** What is the difference between on-policy and off-policy learning? Classify REINFORCE, PPO, and SAC accordingly and explain why.

**Explanation:** The distinction between on-policy and off-policy learning lies in the source of the data used for updates.

**On-Policy Learning:** In on-policy learning, the agent learns from data generated by the *exact same policy* that it is trying to improve. The policy used to collect data (the "behavior policy") is the same as the policy being learned (the "target policy"). This means that after every policy update, the old data is discarded, and new data must be collected with the new policy.

- **Advantage:** More stable, as the updates are directly relevant to the current policy's performance.
- **Disadvantage:** Very sample inefficient, as data is used only once.

**Off-Policy Learning:** In off-policy learning, the agent learns from data generated by a *different* policy than the one it is trying to improve. The behavior policy can be an older version of the target policy, a more exploratory policy, or even data from human demonstrations. This is typically achieved by storing experience in a large replay buffer and sampling from it to perform updates.

- **Advantage:** Highly sample efficient, as data can be reused many times for updates.
- **Disadvantage:** Can be less stable and requires correction techniques (like importance sampling) to account for the distribution mismatch between the behavior and target policies.

#### Classification:

- **REINFORCE: On-policy.** It collects a trajectory with the current policy  $\pi_\theta$ , uses it for a single update, and then discards it.
- **PPO: On-policy.** Like REINFORCE, it collects a batch of data with the current policy. While it may perform multiple gradient steps on this same batch of data, the data is still considered "on-policy" and is discarded before the next round of data collection with the newly updated policy.
- **SAC: Off-policy.** It uses a large experience replay buffer. The agent continuously adds new experience to the buffer and samples random mini-batches of past transitions to update its networks. The data in the buffer comes from many older versions of the policy, making the behavior policy (a mixture of past policies) different from the current target policy being learned.

**Question 12:** Explain the role of the clipped surrogate objective in PPO and how it prevents large, destabilizing policy updates.

**Explanation:** The clipped surrogate objective is the core innovation of PPO, designed to approximate the trust region constraint of TRPO in a simpler, first-order way. The standard policy gradient objective using importance sampling is:

$$L(\theta) = \mathbb{E}_t [r_t(\theta) A_t] \quad \text{where } r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}$$

The problem is that if  $r_t(\theta)$  becomes very large, this can lead to a huge update. PPO introduces a clipped version:

$$L^{\text{CLIP}}(\theta) = \mathbb{E}_t [\min(r_t(\theta) A_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) A_t)]$$

Let's analyze how this works in two cases, based on the sign of the advantage  $A_t$ :

**Case 1: Advantage  $A_t$  is positive.** This means the action was better than average, and we want to increase its probability. The objective becomes:

$$\min(r_t(\theta)A_t, (1 + \epsilon)A_t)$$

The policy update will try to increase  $r_t(\theta)$  to maximize this objective. However, once  $r_t(\theta)$  exceeds  $1 + \epsilon$ , the 'min' function will be determined by the second term,  $(1 + \epsilon)A_t$ . The objective "flattens out," and there is no more incentive to increase the policy ratio further. This effectively "clips" the potential update, preventing the new policy from moving too far from the old one.

**Case 2: Advantage  $A_t$  is negative.** This means the action was worse than average, and we want to decrease its probability. The objective becomes:

$$\min(r_t(\theta)A_t, (1 - \epsilon)A_t)$$

Since  $A_t$  is negative, to maximize the objective (i.e., make it less negative), the update will try to decrease  $r_t(\theta)$ . However, if the update tries to decrease  $r_t(\theta)$  below  $1 - \epsilon$ , the 'min' function will be determined by the first term,  $r_t(\theta)A_t$ . The clipping term  $(1 - \epsilon)A_t$  acts as a lower bound on the objective, again discouraging an overly large change in the policy ratio.

In both cases, the clipping mechanism removes the incentive for the policy update to push the ratio  $r_t(\theta)$  outside the range  $[1 - \epsilon, 1 + \epsilon]$ , thus creating a soft trust region and ensuring more stable learning.

**Question 13:** What is the Policy Improvement Theorem and why is it the bedrock of policy iteration?

**Explanation:** The Policy Improvement Theorem is a fundamental result in reinforcement learning that provides the theoretical guarantee for the "improvement" step in the policy iteration framework.

**The Theorem:** It states that for any two deterministic policies  $\pi$  and  $\pi'$ , if for all states  $s$ , the new policy  $\pi'$  chooses an action that is considered better than or equal to the value of the old policy  $\pi$  at that state, then the new policy  $\pi'$  is guaranteed to be better than or equal to the old policy everywhere.

Formally, if for all  $s \in S$ :

$$Q^\pi(s, \pi'(s)) \geq V^\pi(s)$$

Then it follows that for all  $s \in S$ :

$$V^{\pi'}(s) \geq V^\pi(s)$$

This means the value function of the new policy is monotonically non-decreasing compared to the old one.

**Why it is the Bedrock of Policy Iteration:** Policy Iteration consists of two alternating steps: policy evaluation and policy improvement.

1. In **policy evaluation**, we compute the value function  $V^\pi$  for the current policy  $\pi$ .
2. In **policy improvement**, we create a new policy  $\pi'$  that is "greedy" with respect to  $V^\pi$ . This means for each state  $s$ , the new policy chooses the action  $a$  that maximizes the action-value function:  $\pi'(s) = \arg \max_a Q^\pi(s, a)$ .



The Policy Improvement Theorem is what guarantees this second step actually works. By definition, a greedy policy satisfies the condition  $Q^\pi(s, \pi'(s)) \geq V^\pi(s)$ . Therefore, the theorem guarantees that the new greedy policy  $\pi'$  will be an improvement over (or at least as good as) the original policy  $\pi$ .

This guarantee is the engine of policy iteration. It ensures that by repeatedly evaluating and then improving the policy, we are on a path of monotonic improvement that will eventually converge to the optimal policy,  $\pi^*$ , at which point no further improvements can be made. Without this theorem, there would be no guarantee that the "improvement" step actually leads to a better policy.

**Question 14:** Explain the concept of causality in the context of the REINFORCE algorithm and how using the "reward-to-go" addresses it.

**Explanation:** **Causality** in reinforcement learning is the principle that an action taken at a specific time step,  $t$ , can only influence rewards that occur at future time steps ( $t+1, t+2, \dots$ ) and cannot affect rewards that have already been received (at times  $t, t-1, \dots$ ).

The most basic form of the REINFORCE algorithm violates this principle in its gradient estimator. The initial formulation of the policy gradient is:

$$\nabla_\theta J(\theta) =_{\tau \sim \pi_\theta} \left[ \left( \sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_t | s_t) \right) R(\tau) \right]$$

Here,  $R(\tau) = \sum_{k=0}^T \gamma^k r_{k+1}$  is the total discounted return for the entire episode. When this is distributed, the update term for an action  $a_t$  at time  $t$  is scaled by the full return  $R(\tau)$ . This means the update for  $a_t$  is influenced by rewards  $r_1, r_2, \dots, r_t$  that happened *before* the action was even taken. This makes no logical sense and introduces unnecessary variance into the gradient estimate, as the past rewards are just random noise with respect to the decision at time  $t$ .

The **reward-to-go** estimator corrects this by respecting causality. Instead of scaling the score function  $\nabla_\theta \log \pi_\theta(a_t | s_t)$  by the total return for the whole episode, it scales it by the sum of rewards from that time step onward:

$$\hat{Q}^\pi(s_t, a_t) = \sum_{t'=t}^T \gamma^{t'-t} r_{t'+1}$$

The policy gradient estimator then becomes:

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^T \nabla_\theta \log \pi_\theta(a_{i,t} | s_{i,t}) \hat{Q}_{i,t}^\pi$$

By using the reward-to-go, the update for the action  $a_t$  is now only influenced by the consequences of that action (i.e., future rewards). This does not change the expected value of the gradient (it remains unbiased), but it significantly reduces the variance of the estimate because it removes the irrelevant noise from past rewards. This leads to a more stable and efficient learning signal.

**Question 15:** What is automatic temperature tuning in SAC and what problem does it solve?



**Explanation:** **Automatic temperature tuning** is a feature in modern implementations of Soft Actor-Critic (SAC) that allows the algorithm to learn the optimal value for the temperature parameter,  $\alpha$ , instead of requiring the user to set it as a fixed hyperparameter.

**The Problem it Solves:** The temperature parameter  $\alpha$  is critically important in SAC's maximum entropy objective:

$$J(\pi) = [r(s, a) + \alpha \mathcal{H}(\pi(\cdot|s))]$$

It controls the trade-off between maximizing the reward (exploitation) and maximizing the policy's entropy (exploration).

- If  $\alpha$  is **too high**, the agent will prioritize exploration too much, acting almost randomly and failing to learn the task effectively.
- If  $\alpha$  is **too low**, the agent will ignore the entropy bonus, behave greedily like a standard RL agent, and lose the benefits of exploration and robustness that SAC is designed to provide.

Finding the "just right" value for  $\alpha$  can be difficult and task-dependent, often requiring tedious manual tuning. Furthermore, the optimal balance between exploration and exploitation might change over the course of training.

**How Automatic Tuning Works:** Instead of fixing  $\alpha$ , it is treated as a trainable parameter. A new objective function is defined for  $\alpha$  itself, with the goal of constraining the policy's average entropy to be close to a target value,  $\bar{\mathcal{H}}$ . The target entropy  $\bar{\mathcal{H}}$  is a simple hyperparameter, often set to a heuristic like the negative of the action space dimension ( $-\dim(\mathcal{A})$ ).

The loss function for  $\alpha$  is:

$$J(\alpha) =_{a_t \sim \pi_t} [-\alpha(\log \pi_t(a_t|s_t) + \bar{\mathcal{H}})]$$

This loss is then used to update  $\alpha$  via gradient descent:

- If the policy's current entropy ( $-\log \pi_t$ ) is **higher** than the target entropy  $\bar{\mathcal{H}}$ , the term in the parentheses is positive. The gradient will push  $\alpha$  to decrease, reducing the weight of the entropy bonus and encouraging more exploitation.
- If the policy's current entropy is **lower** than the target entropy, the term is negative. The gradient will push  $\alpha$  to increase, strengthening the entropy bonus and encouraging more exploration.

By dynamically adjusting  $\alpha$  in this way, the algorithm automatically manages the exploration-exploitation trade-off throughout training, removing the need for manual tuning and often leading to better and more stable performance.

**Question 16:** Explain the intellectual lineage from REINFORCE to PPO, highlighting how each step addresses a key limitation of its predecessor.

**Explanation:** The intellectual lineage from REINFORCE to PPO is a clear story of systematically identifying and solving key theoretical and practical problems in on-policy learning.

### 1. REINFORCE (Monte Carlo Policy Gradient):

- **Idea:** The most direct implementation of the Policy Gradient Theorem. It provides a theoretically sound, unbiased estimate of the policy gradient.
- **Limitation: Extremely high variance.** The Monte Carlo return used as the learning signal is very noisy, leading to unstable updates and very poor sample efficiency.

## 2. REINFORCE with Baseline (Advantage Actor-Critic):

- **Idea:** Addresses the high variance of REINFORCE by introducing a state-dependent baseline, optimally the state-value function  $V^\pi(s)$ . This centers the learning signal, creating the Advantage Function  $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s, a)$ . This is the core idea behind most Actor-Critic methods.
- **Limitation:** While variance is reduced, a new problem emerges from the theory: **instability from large updates**. The on-policy approximation used to estimate the gradient is only valid for small policy changes. Large updates can lead to policy collapse.

## 3. TRPO (Trust Region Policy Optimization):

- **Idea:** Directly tackles the instability problem. It formalizes the need for small updates by treating the policy update as a constrained optimization problem. It maximizes the performance objective subject to a hard constraint on how far the new policy can diverge from the old one, measured by KL-divergence. This guarantees monotonic policy improvement.
- **Limitation: High computational complexity.** TRPO is a second-order method that is difficult to implement and computationally expensive, making it impractical for many applications.

## 4. PPO (Proximal Policy Optimization):

- **Idea:** Serves as a practical, first-order approximation of TRPO. It achieves the same goal of stable updates without the complexity. Instead of a hard constraint, it uses a novel **clipped surrogate objective function**.
- **Contribution:** This clipping mechanism effectively creates a soft trust region, discouraging policy updates that are too large. PPO delivers the stability and monotonic improvement guarantees of TRPO in a much simpler, more efficient, and easier-to-implement package, making it the default on-policy algorithm in deep RL today.

This progression shows a clear path: from a sound but impractical idea (REINFORCE), to a more stable version (Advantage Actor-Critic), to a theoretically robust but complex solution for stability (TRPO), and finally to a practical and effective algorithm that captures the best of all worlds (PPO).

**Question 17:** Why is SAC considered more sample-efficient than on-policy methods like PPO?

**Explanation:** SAC (Soft Actor-Critic) is generally much more sample-efficient than on-policy algorithms like PPO primarily because it is an **off-policy** algorithm that uses an **experience replay buffer**.

Here's a breakdown of the reasons:

### 1. Data Re-use (The Main Reason):

- **PPO (On-Policy):** In on-policy learning, data is collected using the current policy  $\pi_\theta$ . After a small number of gradient updates on this data, the policy changes to  $\pi_{\theta'}$ . The collected data is now "stale" and must be discarded. The agent has to go back to the environment to collect new data with  $\pi_{\theta'}$ . Each sample from the environment is used for only a handful of updates.
- **SAC (Off-Policy):** SAC stores all of its experience—transitions of  $(s, a, r, s')$ —in a large data structure called a replay buffer. For each update step, it samples a random mini-batch of transitions from this buffer. This means that a single transition can be re-used for many, many gradient updates as the policy evolves. This ability to "recycle" experience drastically reduces the number of interactions the agent needs with the real environment to learn effectively.

### 2. Decoupling of Exploration and Learning:

- In PPO, the agent must actively explore while learning, and the data it collects is immediately used.
- In SAC, the data collection (adding to the buffer) is decoupled from the learning (sampling from the buffer). This allows the agent to learn from a more diverse set of experiences, including highly successful transitions from the distant past, which can speed up learning.

### 3. Improved Exploration from Maximum Entropy:

- While not directly related to being off-policy, SAC's maximum entropy objective encourages more systematic and robust exploration. This can lead the agent to discover rewarding parts of the state space more quickly, making the data it collects more valuable and thus improving efficiency.

In summary, the core reason for SAC's superior sample efficiency is its ability to learn from each piece of experience multiple times via the replay buffer, a hallmark of off-policy learning that on-policy methods like PPO cannot leverage.

**Question 18:** What is the reparameterization trick and why is it necessary for training the actor in SAC?

**Explanation:** The **reparameterization trick** is a technique used to compute a low-variance gradient for a stochastic node in a computation graph. It is essential for training the actor (the policy network) in modern actor-critic methods like SAC.

**The Problem:** The actor in SAC is a stochastic policy,  $\pi_\phi(a|s)$ , which means it outputs a probability distribution (e.g., a Gaussian with mean  $\mu_\phi(s)$  and standard deviation  $\sigma_\phi(s)$ ). To train the actor, we need to update its parameters  $\phi$  to produce actions that lead to higher soft Q-values. This requires calculating the gradient of the Q-function with respect to the policy parameters  $\phi$ .

The process looks like this: 1. Sample an action:  $a \sim \pi_\phi(\cdot|s)$  2. Evaluate its Q-value:  $Q_\theta(s, a)$  3. Compute the gradient:  $\nabla_\phi Q_\theta(s, a)$

The issue is step 1. Sampling is a random, non-differentiable operation. You cannot backpropagate a gradient through a sampling node. The REINFORCE algorithm gets around this using the log-derivative trick, but that method is known to have high variance.

**The Solution: Reparameterization Trick** The trick reformulates the sampling process to isolate the randomness. Instead of sampling the action directly from the policy's output distribution, we sample a random noise vector from a simple, fixed distribution (like a standard normal distribution  $\mathcal{N}(0, I)$ ) and then transform this noise using a deterministic function parameterized by the policy network.

For a Gaussian policy, the process becomes: 1. Sample random noise:  $\epsilon \sim \mathcal{N}(0, I)$  2. Compute the action deterministically:  $a = \mu_\phi(s) + \sigma_\phi(s) \odot \epsilon$

Now, the action  $a$  is a deterministic function of the network's outputs  $(\mu_\phi, \sigma_\phi)$  and the external noise vector  $\epsilon$ . The stochasticity has been moved outside the network.

**Why it is Necessary for SAC:** This reformulation allows us to compute the policy gradient easily and with low variance using standard backpropagation. The gradient  $\nabla_\phi Q_\theta(s, a)$  can now flow through the Q-network, back through the deterministic function for  $a$ , and into the policy network parameters  $\phi$  via the chain rule.

This provides a much more stable and efficient gradient estimate for the policy update compared to the high-variance score function estimator used in REINFORCE. This stability is crucial for the performance of off-policy actor-critic algorithms like SAC.

**Question 19:** Explain the two main processes of Soft Policy Iteration (Evaluation and Improvement) and how they guarantee convergence.

**Explanation:** Soft Policy Iteration is the theoretical framework that guarantees the convergence of Soft Actor-Critic (SAC). It adapts the classic Generalized Policy Iteration (GPI) framework to the maximum entropy objective. It consists of two alternating steps: Soft Policy Evaluation and Soft Policy Improvement.

**1. Soft Policy Evaluation (Lemma 1):** This step aims to compute the "soft" Q-function for a given policy  $\pi$ . The soft Q-function,  $Q^\pi$ , represents the expected future return plus the expected future entropy. It is the fixed point of the **soft Bellman backup operator**,  $\mathcal{T}^\pi$ :

$$\mathcal{T}^\pi Q(s, a) \triangleq r(s, a) + \gamma_{s', a' \sim \pi} [Q(s', a') - \alpha \log \pi(a'|s')]$$

Lemma 1 proves that this operator  $\mathcal{T}^\pi$  is a **contraction mapping**. A key property of contraction mappings is that if you apply them repeatedly to any initial Q-function, they are guaranteed to converge to a unique fixed point. In this case, the fixed point is the true soft Q-function,  $Q^\pi$ . This step answers the question: "How good is my current policy under the maximum entropy objective?"

**2. Soft Policy Improvement (Lemma 2):** This step aims to find a new, better policy based on the soft Q-function computed in the evaluation step. Instead of a simple greedy update, the new policy  $\pi_{\text{new}}$  is chosen to be the one that minimizes the KL-divergence to the Boltzmann distribution (the exponential) of the old policy's soft Q-function:

$$\pi_{\text{new}} = \arg \min_{\pi' \in \Pi} \left( \pi'(\cdot|s) \left\| \frac{\exp(\frac{1}{\alpha} Q^{\pi_{\text{old}}}(s, \cdot))}{Z(s)} \right\| \right)$$

Lemma 2 proves that this update guarantees **monotonic improvement**. That is, the soft Q-value of the new policy is guaranteed to be greater than or equal to the soft Q-value of the old policy for all state-action pairs:  $Q^{\pi_{\text{new}}}(s, a) \geq Q^{\pi_{\text{old}}}(s, a)$ . This step answers the question: "How can I improve my policy to better balance rewards and entropy?"

**Guarantee of Convergence (Theorem 1):** The overall convergence of Soft Policy Iteration is guaranteed by the interaction of these two steps.

1. Lemma 2 ensures that each policy improvement step produces a new policy whose soft Q-function is monotonically non-decreasing.
2. The soft Q-values are bounded above (assuming bounded rewards).
3. A sequence that is both monotonically non-decreasing and bounded above is guaranteed to converge.

Therefore, by repeatedly applying soft evaluation and soft improvement, the sequence of soft Q-functions converges to the optimal soft Q-function,  $Q^*$ , and the policy converges to the optimal soft policy,  $\pi^*$ , within the considered class of policies.

**Question 20:** What is the Clipped Double-Q Learning trick in SAC, and what specific problem in Q-learning does it address?

**Explanation:** **Clipped Double-Q Learning** is a technique used in modern actor-critic algorithms like SAC and TD3 to improve the stability of learning the critic (the Q-function). It involves two key ideas: using two Q-networks (twin critics) and taking the minimum of their predictions for the target value calculation.

**The Problem it Addresses: Q-Value Overestimation** Standard Q-learning and many deep Q-learning algorithms (like DDPG) suffer from a systematic **overestimation bias**. The update for the Q-function involves a maximization step over the next state's Q-values when calculating the target (e.g.,  $y = r + \gamma \max_{a'} Q(s', a')$ ).

Because the Q-values are just estimates and contain noise, this maximization is performed over noisy values. The maximum of a set of noisy values is likely to be higher than the true maximum value. This consistent overestimation can propagate through the learning process, leading the algorithm to learn an incorrect Q-function and, consequently, a suboptimal policy that exploits non-existent high values.

**How Clipped Double-Q Learning Works in SAC:** SAC implements a solution to this problem with three components:

1. **Twin Critics:** Instead of one Q-network, SAC learns two independent Q-networks,  $Q_{\theta_1}$  and  $Q_{\theta_2}$ , simultaneously. They are trained on the same data but their independent initializations and stochastic updates cause them to learn slightly different values.
2. **Twin Target Networks:** Each critic has its own corresponding target network,  $Q_{\theta_{1,\text{targ}}}$  and  $Q_{\theta_{2,\text{targ}}}$ , for stable learning.
3. **Clipped (Minimum) Target Value:** This is the core of the trick. When calculating the target value  $y$  for the Bellman update, SAC computes the next-state value using both target networks and takes the **minimum** of the two.

$$y = r_t + \gamma \left( \min_{j=1,2} Q_{\theta_{j,\text{targ}}}(s_{t+1}, a_{t+1}) - \alpha \log \pi_{\phi}(a_{t+1}|s_{t+1}) \right)$$

By taking the minimum, the algorithm creates a more conservative, pessimistic estimate of the future Q-value. It makes it less likely that a single spuriously high Q-value estimate from one network will be used to update the other, thus mitigating the overestimation bias. This leads to more accurate value estimates, more stable training, and ultimately better final performance.

**Question 21:** Explain the role of the Advantage Function as a refined learning signal compared to the total return or reward-to-go.

**Explanation:** The Advantage Function,  $A^\pi(s, a)$ , represents a principled refinement of the learning signal in policy gradient methods, designed to create more stable and efficient learning by reducing variance.

Let's trace the evolution of the learning signal: **1. Total Return ( $R(\tau)$ ):** This is the most basic signal, used in the original formulation of REINFORCE. It scales the policy gradient by the total reward of the entire episode.

- **Problem:** It has extremely high variance and ignores causality (an action at time  $t$  is affected by rewards from before  $t$ ).

**2. Reward-to-Go ( $\hat{Q}^\pi(s_t, a_t)$ ):** This is the first refinement. It scales the gradient for an action at time  $t$  only by the rewards that come after it.

- **Improvement:** It respects causality, which reduces some variance.
- **Problem:** It still has high variance. The absolute value of the return can vary wildly depending on the episode, even if the action was good relative to other actions in that state. For example, in a game where scores range from 1,000 to 1,100, a return of 1,005 is a "good" outcome and will be reinforced, but it's much worse than an outcome of 1,090. The signal isn't very discriminative.

**3. Advantage Function ( $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$ ):** This is the final and most refined signal. It subtracts the state-value function,  $V^\pi(s)$ , from the action-value function,  $Q^\pi(s, a)$ .

- **Improvement:** The Advantage Function answers the question: "**How much better was this specific action compared to the average action I usually take in this state?**"
- If  $A^\pi(s, a) > 0$ , the action was better than average, and its probability should be increased.
- If  $A^\pi(s, a) < 0$ , the action was worse than average, and its probability should be decreased.
- If  $A^\pi(s, a) \approx 0$ , the action was average, and there's no strong reason to change its probability.
- **Benefit:** This relative signal has much lower variance because it removes the shared, baseline component of the return ( $V^\pi(s)$ ) and focuses only on the marginal contribution of the specific action. It is invariant to reward shifting (adding a constant to all rewards doesn't change the advantage) and provides a much more stable and informative learning signal, leading to faster and more reliable convergence.

**Question 22:** What is the Performance Difference Lemma and what is its significance in connecting policy gradients to policy iteration?

**Explanation:** The **Performance Difference Lemma** is a crucial theoretical result in modern policy gradient theory. It provides an exact expression for the difference in performance between two policies,  $\pi_{\theta'}$  and  $\pi_{\theta}$ .

**The Lemma:** It states that the change in the objective function when switching from policy  $\pi_{\theta}$  to  $\pi_{\theta'}$  is equal to the expected sum of the advantages of the new policy's actions, where the advantages are calculated with respect to the old policy  $\pi_{\theta}$ .

$$J(\theta') - J(\theta) =_{\tau \sim \pi_{\theta'}} \left[ \sum_{t=0}^{\infty} \gamma^t A^{\pi_{\theta}}(s_t, a_t) \right]$$

Note that the expectation is taken over trajectories generated by the *new* policy,  $\pi_{\theta'}$ .

**Significance and Connection to Policy Iteration:** The lemma's significance lies in how it bridges the gap between two seemingly different families of RL algorithms: policy gradient methods and policy iteration.

1. **Justifies the Advantage Function:** The lemma shows that the quantity we want to maximize to improve our policy is precisely the expected advantage. This provides a strong theoretical justification for using the advantage function as the learning signal in policy gradient methods. The policy gradient update,  $\nabla_{\theta} J(\theta) \propto [\nabla_{\theta} \log \pi_{\theta}(a|s) A^{\pi_{\theta}}(s, a)]$ , is directly taking a step in the direction that maximizes this performance difference.
2. **Frames Policy Gradient as Policy Iteration:** The lemma allows us to see policy gradient methods as a form of Generalized Policy Iteration (GPI):
  - **Policy Evaluation:** In this step, we estimate the advantage function  $A^{\pi_{\theta}}(s, a)$  for the current policy. This is analogous to the evaluation step in GPI.
  - **Policy Improvement:** The policy gradient update then uses this advantage estimate to find a better policy. This is analogous to the improvement step in GPI.

The lemma makes this connection formal. It shows that the policy gradient update is directly optimizing the very quantity that the Policy Improvement Theorem in classic GPI is based on. It reveals that policy gradient methods are essentially performing an "infinitesimal" version of policy iteration, where each update is a small step towards a better policy rather than a large, discrete jump.

**Question 23:** Why is importance sampling theoretically sound for correcting distribution mismatch but often impractical for on-policy methods?

**Explanation:** **Importance Sampling (IS)** is a statistical technique that allows us to estimate the expectation of a function under a target distribution ( $p$ ) using samples drawn from a different behavior distribution ( $q$ ). The formula is:

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

The ratio  $\frac{p(x)}{q(x)}$  is the importance weight.



**Why it is Theoretically Sound:** In the context of policy gradients, we want to evaluate the performance of a new policy  $\pi_{\theta'}$  using data from an old policy  $\pi_{\theta}$ . Importance sampling provides the exact mathematical correction for this distribution mismatch. The corrected objective for the performance difference would involve multiplying the advantage at each step by the ratio of action probabilities:  $\frac{\pi_{\theta'}(a_t|s_t)}{\pi_{\theta}(a_t|s_t)}$ . This correctly re-weights the samples to give an unbiased estimate of the objective under the new policy.

**Why it is Often Impractical:** The practicality of importance sampling breaks down due to the problem of **high variance**, especially when correcting for the full trajectory distribution.

A full correction for the state distribution mismatch requires considering the probability of the entire trajectory. The importance weight for a trajectory  $\tau$  is:

$$\frac{p_{\theta'}(\tau)}{p_{\theta}(\tau)} = \frac{\prod_{t=0}^T \pi_{\theta'}(a_t|s_t) p(s_{t+1}|s_t, a_t)}{\prod_{t=0}^T \pi_{\theta}(a_t|s_t) p(s_{t+1}|s_t, a_t)} = \prod_{t=0}^T \frac{\pi_{\theta'}(a_t|s_t)}{\pi_{\theta}(a_t|s_t)}$$

This weight is a product of ratios over the entire episode horizon  $T$ .

- **Exponential Growth of Variance:** If the policies  $\pi_{\theta'}$  and  $\pi_{\theta}$  are even slightly different at each step, the product of these ratios can either explode to a very large number or vanish to zero.
- **Unreliable Estimates:** As a result, the variance of the importance-weighted estimator grows exponentially with the length of the episode. In practice, this means the estimate will be dominated by a very small number of trajectories that happen to get an enormous weight, while most others get a weight close to zero. The resulting gradient estimate is extremely noisy and unreliable, making learning impossible.

Because of this exponential variance problem, full trajectory importance sampling is rarely used. On-policy methods like PPO use a single-step importance ratio and then clip it to keep the variance under control, but this is an approximation, not a full correction. Off-policy methods like SAC avoid this issue by using a replay buffer and a different theoretical framework (Soft Policy Iteration) that doesn't rely on these high-variance trajectory-wise corrections.

**Question 24:** Describe the complete update process for one step in the SAC algorithm, explaining the calculation of all three loss functions ( $J_Q$ ,  $J_{\pi}$ , and  $J_{\alpha}$ ).

**Explanation:** A single update step in the Soft Actor-Critic (SAC) algorithm involves sampling a mini-batch of transitions from the replay buffer and updating the critic, actor, and temperature parameters by minimizing their respective loss functions.

Let's assume we have sampled a mini-batch of transitions  $(s, a, r, s', d)$  from the replay buffer  $\mathcal{D}$ , where  $d$  is a done flag.

**Step 1: Critic Update (Minimizing  $J_Q$ )** The goal is to make the Q-networks' predictions match the soft Bellman target.

#### 1. Compute the Target Value ( $y$ ):

- Sample a next action  $a'$  from the *current* policy:  $a' \sim \pi_{\phi}(\cdot|s')$ .
- Get the log-probability of this next action:  $\log \pi_{\phi}(a'|s')$ .



- Evaluate the next state's value using the *target* Q-networks and take the minimum:  $Q_{\min\_targ} = \min_{j=1,2} Q_{\theta_j, targ}(s', a')$ .
- Combine these to form the target  $y$ :

$$y = r + \gamma(1 - d)(Q_{\min\_targ} - \alpha \log \pi_\phi(a'|s'))$$

The  $(1 - d)$  term ensures the future value is zero if the state is terminal.

2. **Compute the Q-Loss:** The loss for each critic is the mean squared error between its prediction for the current state-action pair  $(s, a)$  and the target  $y$ .

$$J_{Q_1} =_{(s,a) \sim \mathcal{D}} \left[ \frac{1}{2} (Q_{\theta_1}(s, a) - y)^2 \right]$$

$$J_{Q_2} =_{(s,a) \sim \mathcal{D}} \left[ \frac{1}{2} (Q_{\theta_2}(s, a) - y)^2 \right]$$

The total Q-loss is  $J_Q = J_{Q_1} + J_{Q_2}$ .

3. **Update Critics:** Perform a gradient descent step on the critic parameters  $\theta_1$  and  $\theta_2$  using this loss.

**Step 2: Actor Update (Minimizing  $J_\pi$ )** The goal is to update the policy to output actions that maximize the soft Q-value.

1. **Compute Policy Actions:** Sample actions for the states  $s$  from the current batch using the reparameterization trick:  $a_{\text{new}} \sim \pi_\phi(\cdot|s)$ .
2. **Compute the Policy Loss:** The loss encourages actions with high Q-values and high entropy.
  - Get the log-probability of these new actions:  $\log \pi_\phi(a_{\text{new}}|s)$ .
  - Evaluate the Q-value of these new actions using one of the critics (e.g., the minimum of the two):  $Q_{\min} = \min_{j=1,2} Q_{\theta_j}(s, a_{\text{new}})$ .
  - The policy loss is:

$$J_\pi =_{s \sim \mathcal{D}} [\alpha \log \pi_\phi(a_{\text{new}}|s) - Q_{\min}]$$

3. **Update Actor:** Perform a gradient ascent step (or descent on  $-J_\pi$ ) on the actor parameters  $\phi$ .

**Step 3: Temperature Update (Minimizing  $J_\alpha$ )** The goal is to adjust  $\alpha$  to keep the policy's entropy near a target value  $\bar{\mathcal{H}}$ .

1. **Compute the Alpha Loss:**

$$J_\alpha =_{s \sim \mathcal{D}} [-\alpha(\log \pi_\phi(a_{\text{new}}|s) + \bar{\mathcal{H}})]$$

Note: We use the same log-probabilities calculated for the policy loss.

2. **Update Alpha:** Perform a gradient descent step on the (log of)  $\alpha$ .

Finally, the target networks are updated using Polyak averaging:

$$\theta_{j,\text{targ}} \leftarrow \rho\theta_{j,\text{targ}} + (1 - \rho)\theta_j$$

This completes one full update cycle.

**Question 25:** What is the difference between a stochastic and a deterministic policy? Why are stochastic policies central to the policy gradient methods discussed?

**Explanation:** The difference between a stochastic and a deterministic policy lies in how they map states to actions.

**Deterministic Policy:** A deterministic policy, denoted  $\mu(s)$ , maps each state directly to a single action. Given a state  $s$ , the action is always the same:  $a = \mu(s)$ . There is no randomness in the action selection process itself. Algorithms like DDPG learn deterministic policies.

**Stochastic Policy:** A stochastic policy, denoted  $\pi(a|s)$ , maps each state to a probability distribution over actions. Given a state  $s$ , the policy outputs the probability of taking each possible action. The agent then samples from this distribution to select its action. This means that in the same state, the agent might take different actions on different visits.

**Why Stochastic Policies are Central to Policy Gradient Methods:** Stochastic policies are fundamental to the policy gradient methods discussed in the report (REINFORCE, PPO, SAC) for several key reasons:

1. **Gradient Calculation:** The Policy Gradient Theorem itself is derived for stochastic policies. The gradient is formulated in terms of  $\log \pi_\theta(a|s)$ , the log-probability of taking an action. This term is only well-defined for a policy that outputs a distribution. For a deterministic policy, the probability of taking the chosen action is 1, and 0 for all others, which creates mathematical difficulties for this formulation.
2. **Exploration:** Stochastic policies have built-in exploration. By nature, they try different actions according to their probabilities. A deterministic policy, on the other hand, has no built-in exploration. To explore, it must have external noise added to its actions, which can be less systematic. The entropy of a stochastic policy, which SAC explicitly maximizes, is a direct measure of this exploration.
3. **Handling Uncertainty and Partial Observability:** In many real-world problems, the state representation may be incomplete or noisy (partially observable). A stochastic policy can be optimal in such cases. For example, in rock-paper-scissors, the optimal policy is to play each action with 1/3 probability. A deterministic policy would be easily exploited. A stochastic policy can better handle this uncertainty by hedging its bets.
4. **Smoother Learning Landscape:** Stochastic policies can lead to a smoother optimization landscape. A small change in the parameters of a stochastic policy leads to a small change in the action probabilities, whereas a small change in a deterministic policy's parameters could cause it to switch abruptly from one action to another, leading to more volatile learning.

While some off-policy methods learn deterministic policies, the entire on-policy policy gradient framework and the maximum entropy framework of SAC are built upon the foundation of learning and optimizing a stochastic policy.

**Question 26:** Explain the concept of a "trust region" and how TRPO's constrained optimization problem formalizes it.

**Explanation:** A **trust region** in optimization is a conceptual region around the current set of parameters within which a simplified model (like a linear or quadratic approximation) of the objective function is considered trustworthy or reliable. When we take an update step, we constrain the step to stay within this region to ensure that the update leads to a real improvement and doesn't "fall off a cliff" where the approximation is no longer valid.

In the context of policy gradients, the "simplified model" is the performance improvement objective calculated using data from the old policy  $\pi_\theta$ :

$$L_\theta(\theta') =_{s,a \sim \pi_\theta} \left[ \frac{\pi_{\theta'}(a|s)}{\pi_\theta(a|s)} A^{\pi_\theta}(s, a) \right]$$

As we know from the state distribution mismatch problem, this approximation  $L_\theta(\theta')$  is only a reliable estimate of the true performance improvement  $J(\theta') - J(\theta)$  if the new policy  $\pi_{\theta'}$  is "close" to the old policy  $\pi_\theta$ . The trust region, therefore, is the set of new policies that are close enough to the old one for this approximation to hold.

**How TRPO Formalizes the Trust Region:** TRPO (Trust Region Policy Optimization) turns this conceptual idea into a formal, constrained optimization problem. At each step, it aims to maximize the approximate objective  $L_\theta(\theta')$  while explicitly forcing the new policy to stay within the trust region.

It defines the "closeness" of policies using the **Kullback-Leibler (KL) divergence**, a measure of the difference between two probability distributions. The TRPO update rule is:

$$\begin{aligned} &\text{maximize}_{\theta'} && L_\theta(\theta') \\ &\text{subject to} && \int_{s \sim \pi_\theta} [\pi_\theta(\cdot|s) \log \pi_{\theta'}(\cdot|s)] \leq \delta \end{aligned}$$

Here:

- The **objective** is to find the policy parameters  $\theta'$  that maximize the estimated performance improvement.
- The **constraint** is the formal definition of the trust region. It requires that the average KL-divergence between the old policy and the new policy, over the states visited by the old policy, must be less than or equal to a small constant  $\delta$ .

By solving this problem, TRPO finds the largest possible improvement step that it can "trust" to be valid. This provides a theoretical guarantee of monotonic policy improvement and prevents the catastrophic policy collapse that can happen with unconstrained updates.

**Question 27:** What is the difference between Monte Carlo and Temporal Difference (TD) learning, and how do they relate to the estimation of value functions in policy gradient methods?

**Explanation:** Monte Carlo (MC) and Temporal Difference (TD) are two fundamental methods in reinforcement learning for estimating value functions from experience. They differ primarily in how and when they update their estimates.

#### Monte Carlo (MC) Learning:

- **Update Rule:** MC methods wait until the **end of an entire episode** to update the value function. The value of a state is estimated as the average of the full, observed returns (rewards-to-go) from all visits to that state.
- **Properties:**
  - **Unbiased:** The estimate is an unbiased sample of the true value function.
  - **High Variance:** The estimate has high variance because the full return depends on many random actions and state transitions that occur over the rest of the episode.
  - **Episodic Tasks Only:** It can only be applied to tasks that have a clear end.

#### Temporal Difference (TD) Learning:

- **Update Rule:** TD methods update the value function after **every single time step**. They do not wait for the end of the episode. They update their current estimate for a state's value,  $V(s_t)$ , based on the observed reward  $r_{t+1}$  and their current estimate of the next state's value,  $V(s_{t+1})$ . This is called **bootstrapping**. The TD(0) update is:

$$V(s_t) \leftarrow V(s_t) + \alpha(r_{t+1} + \gamma V(s_{t+1}) - V(s_t))$$

- **Properties:**
  - **Biased:** The update relies on an existing, likely incorrect estimate ( $V(s_{t+1})$ ), so the update target is biased.
  - **Low Variance:** The update has much lower variance than MC because it only depends on one random action and one state transition.
  - **Continuous Tasks:** It can be used for continuous (non-terminating) tasks.

**Relation to Policy Gradient Methods:** These two learning methods directly correspond to how the advantage function is estimated in different policy gradient algorithms.

- **REINFORCE** is a pure **Monte Carlo** method. It estimates the advantage (or reward-to-go) by running a full episode and calculating the actual, full return. This is why it suffers from high variance.
- **Actor-Critic methods (like A2C, PPO, SAC)** are typically based on **Temporal Difference** learning. The "critic" is a value function approximator that is trained using TD methods. The advantage function is then estimated using the critic's bootstrapped value estimates. For example, a common advantage estimate is the TD error:

$$\hat{A}(s_t, a_t) = r_{t+1} + \gamma V_\omega(s_{t+1}) - V_\omega(s_t)$$

where  $V_\omega$  is the critic network. This use of TD learning is what allows actor-critic methods to be more sample-efficient and have lower variance than REINFORCE.

**Question 28:** Why is policy entropy (and thus the temperature  $\alpha$ ) generally bounded in discrete action spaces but can be unbounded in continuous spaces?

**Explanation:** The boundedness of policy entropy,  $\mathcal{H}(\pi(\cdot|s)) =_{a \sim \pi} [-\log \pi(a|s)]$ , depends on the nature of the action space.

**Discrete Action Spaces:** In a discrete action space with  $K$  possible actions, the policy  $\pi(a|s)$  is a probability vector  $[p_1, p_2, \dots, p_K]$  where  $\sum p_i = 1$ .

- **Maximum Entropy:** The entropy is maximized when the policy is a uniform distribution, i.e.,  $p_i = 1/K$  for all actions. In this case, the entropy is:

$$\mathcal{H}_{\max} = - \sum_{i=1}^K \frac{1}{K} \log \left( \frac{1}{K} \right) = - \log \left( \frac{1}{K} \right) = \log(K)$$

This is a finite, constant value.

- **Minimum Entropy:** The entropy is minimized when the policy is deterministic, i.e., one  $p_i = 1$  and all others are 0. The entropy is  $-\log(1) = 0$ .

Since the entropy is always between 0 and  $\log(K)$ , it is **bounded**. This means the soft Bellman backup operator in SAC is guaranteed to be a contraction, ensuring convergence.

**Continuous Action Spaces:** In a continuous action space, the policy is typically represented by a probability density function (pdf), most commonly a Gaussian distribution  $\mathcal{N}(\mu, \sigma^2)$ . The entropy of a Gaussian distribution is given by:

$$\mathcal{H} = \frac{1}{2} \log(2\pi e \sigma^2)$$

- **Unbounded Nature:** The entropy depends directly on the variance,  $\sigma^2$ . The policy network can learn to make the variance arbitrarily large. As  $\sigma^2 \rightarrow \infty$ , the entropy  $\mathcal{H} \rightarrow \infty$ .
- **Implications for SAC:** Because the entropy can be infinite, the "entropy-augmented reward" in the soft Bellman backup is not guaranteed to be bounded. This technically breaks the proof that the soft Bellman operator is a contraction mapping. In practice, this is not a major issue because:
  1. The policy is incentivized to keep the variance from exploding because doing so would lead to poor performance on the reward part of the objective.
  2. The automatic temperature tuning mechanism in SAC helps to control the policy's entropy, preventing it from growing without bound by adjusting  $\alpha$ .

So, while there is a theoretical gap for continuous spaces, practical implementations of SAC work very well due to these implicit and explicit regularization effects.

**Question 29:** Trace the evolution of the learning signal from REINFORCE to Advantage Actor-Critic, explaining the motivation for each refinement.

**Explanation:** The evolution of the learning signal from REINFORCE to Advantage Actor-Critic is a story of progressively reducing variance to achieve more stable and efficient learning.

**Stage 1: Total Episodic Return (Basic REINFORCE)**

- **Signal:**  $G_0 = \sum_{t=0}^T \gamma^t r_{t+1}$
- **Update Form:**  $\nabla_{\theta} J(\theta) \propto [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \cdot G_0]$
- **Motivation:** This is the most direct, unbiased implementation of the Policy Gradient Theorem.
- **Limitation:** It has extremely high variance and ignores causality. The update for an action at time  $t$  is influenced by rewards that came before it, which makes no sense and adds noise.

### Stage 2: Reward-to-Go (Refined REINFORCE)

- **Signal:**  $G_t = \sum_{k=t}^T \gamma^{k-t} r_{k+1}$
- **Update Form:**  $\nabla_{\theta} J(\theta) \propto [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \cdot G_t]$
- **Motivation:** To fix the causality problem. An action at time  $t$  should only be judged by its consequences (rewards from  $t$  onward).
- **Limitation:** While it reduces variance by removing irrelevant past rewards, the signal still has high variance. The absolute magnitude of  $G_t$  can fluctuate wildly between episodes, making it hard to tell if an action was good in a relative sense.

### Stage 3: Advantage Function (Advantage Actor-Critic)

- **Signal:**  $A^{\pi}(s_t, a_t) = G_t - V^{\pi}(s_t)$
- **Update Form:**  $\nabla_{\theta} J(\theta) \propto [\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \cdot A^{\pi}(s_t, a_t)]$
- **Motivation:** To address the remaining high variance by changing the question from "How good was the total outcome?" to "How much better was this action than the average action in this state?".
- **Contribution:** This is achieved by subtracting a state-dependent baseline, ideally the state-value function  $V^{\pi}(s_t)$ . This centers the learning signal around zero. Positive values correspond to better-than-average actions, and negative values to worse-than-average actions. This relative signal is much more informative and has significantly lower variance, leading to the stable and efficient learning characteristic of modern Actor-Critic methods.

**Question 30:** If policy gradient methods are an "infinitesimal" form of policy iteration, what would a "discrete" or "large step" version look like, and why is it generally not used with function approximation?

**Explanation:** Policy gradient methods are considered "infinitesimal" because they take a small step in the direction of policy improvement. A "discrete" or "large step" version would be a direct implementation of classic Policy Iteration with function approximation.

**What a "Large Step" Version Would Look Like:** The process would follow the classic GPI loop:

1. **Policy Evaluation:** Given the current policy network  $\pi_\theta$ , train a Q-function network  $Q_\omega(s, a)$  until it accurately approximates the true action-value function  $Q^{\pi_\theta}(s, a)$ . This would involve many updates to the critic network until it converges for the fixed policy.
2. **Policy Improvement:** Create a new, completely updated policy  $\pi_{\theta'}$  by making it greedy with respect to the learned Q-function. For every state  $s$ , the new policy would be:

$$\pi_{\theta'}(s) = \arg \max_a Q_\omega(s, a)$$

This would involve training the policy network  $\pi_{\theta'}$  to output the argmax action, which could be a large, discrete change in the policy's behavior.

This loop would then be repeated.

**Why This is Generally Not Used with Function Approximation:** This discrete, large-step approach, while sound in tabular settings, is highly unstable and often fails when using non-linear function approximators like deep neural networks. There are two main reasons:

**1. Instability of Greedy Updates:** The policy improvement step,  $\pi_{\theta'}(s) = \arg \max_a Q_\omega(s, a)$ , is a very aggressive update. The Q-function  $Q_\omega$  is only an approximation of the true value function and inevitably contains errors. When we take a hard 'argmax' over these noisy, imperfect values, we are engaging in "maximization over a noisy set," which can amplify the errors. The new policy might exploit these errors, leading it to a region of the state space where its performance is actually terrible, causing the learning process to diverge.

**2. Moving Target Problem Distribution Shift:** The two steps are fundamentally at odds when using function approximation.

- The policy evaluation step tries to learn the value of a fixed policy  $\pi_\theta$ .
- The policy improvement step makes a large change to the policy, creating  $\pi_{\theta'}$ .

This large change means the state distribution shifts dramatically (the state distribution mismatch problem in its most extreme form). The Q-function  $Q_\omega$  that was painstakingly learned for  $\pi_\theta$  is now largely irrelevant for evaluating the new policy  $\pi_{\theta'}$ . The process becomes unstable because the critic is always trying to learn a "moving target."

The infinitesimal approach of policy gradients is a solution to this instability. By taking small steps, the policy changes slowly, the state distribution shifts gradually, and the critic can keep up with the changes in the actor, leading to a much more stable, co-adapted learning process.