# A Foundational Monograph on Modern Reinforcement Learning

## From Function Approximation to Policy Optimization

by Taha Majlesi

July 17, 2025

# Contents

# Part I

# The Challenge of Scale and Value-Based Approximation

# Chapter 1

# The Imperative for Function Approximation

The central objective of reinforcement learning (RL) is to develop agents that can learn optimal behaviors through interaction with an environment. In its classical formulation, this learning process is often managed using tabular methods, where the agent maintains an explicit record of the value associated with every possible state or state-action pair. While theoretically robust, this approach confronts a practical and insurmountable barrier when applied to problems of realistic complexity: the **curse of dimensionality**. This section establishes the fundamental challenge posed by large-scale environments and introduces the paradigm of function approximation as the necessary solution. It builds the conceptual and mathematical framework for value-based approximation, culminating in an in-depth analysis of the Deep Q-Network (DQN), a landmark algorithm that made large-scale value-based learning viable.

## 1.1 The Curse of Dimensionality in Tabular RL

Traditional reinforcement learning algorithms, such as Q-learning or SARSA, were originally conceived for environments with a finite and manageable number of states and actions. In this context, the learning problem can be solved by maintaining a lookup table. For a state-value function, this table would have an entry, $v(s)$, for every state $s$ in the environment's state space $S$. Similarly, for an action-value function, the table would store an entry, $q(s, a)$, for every state-action pair $(s, a)$. The agent's learning process involves iteratively updating these table entries based on the rewards it receives from the environment, eventually converging to the true optimal values.

This tabular approach, however, collapses under the weight of large-scale problems. The number of states can grow exponentially with the number of variables describing the environment, a phenomenon known as the **curse of dimensionality**. Consider the state spaces of now-classic RL benchmarks:

- The board game Backgammon has an estimated $10^{20}$ states.

- The game of Go has a staggering $10^{170}$ states.

- Problems with continuous state spaces, such as controlling a robotic arm, have an infinite number of unique states.

- In computer vision, a 128x128 pixel RGB image constitutes a state space with approximately $256^{128 \times 128 \times 3} \approx 10^{74000}$ possible states.

The practical consequences of this state space explosion are twofold. First, the memory required to store the lookup table becomes prohibitive. Second, and more critically, the time required to learn the value of each state individually becomes astronomical. This intractability is the primary causal driver for the development of **function approximation**. The goal must shift from *storing* an exact value for every state to *approximating* the underlying value function.

## 1.2 The Core Idea: Generalizing with Parameterized Functions

The solution to the curse of dimensionality is to approximate the value function with a parameterized function. This function, denoted by a vector of parameters or weights $\mathbf{w}$ (often written as $\theta$), captures the underlying relationship between states and their values, allowing the agent to generalize from states it has seen to those it has not.

This concept can be formalized as follows:

- **State-value function approximation:** We learn a function $\hat{v}(s; \mathbf{w})$ that aims to be a close approximation of the true state-value function, $v_\pi(s)$.

$$\hat{v}(s; \mathbf{w}) \approx v_\pi(s) \tag{1.1}$$

- **Action-value function approximation:** Similarly, we learn a function $\hat{Q}(s, a; \mathbf{w})$ to approximate the true action-value function, $q_\pi(s, a)$.

$$\hat{Q}(s, a; \mathbf{w}) \approx q_\pi(s, a) \tag{1.2}$$

The parameter vector $\mathbf{w}$ has a dimensionality that is far smaller than the number of states, $|S|$, which resolves the memory issue. More importantly, by learning a functional form, a single update to the weights $\mathbf{w}$ based on an experience in one state will influence the predicted values for many other states, enabling generalization.

## 1.3 The Objective: Minimizing Mean Squared Value Error (MSVE)

To learn the optimal parameters $\mathbf{w}$, we define a loss function that quantifies the discrepancy between our approximation and the true value function. The most common choice is the **Mean Squared Value Error (MSVE)**. The objective function, $J(\mathbf{w})$, is the MSVE weighted by the stationary distribution of states under the policy $\pi$, denoted as $\mu_\pi(s)$. This distribution represents the fraction of time the agent spends in state $s$ in the long run.

$$J(\mathbf{w}) = \mathbb{E}_{s \sim \mu_\pi} \left[ (v_\pi(s) - \hat{v}(s; \mathbf{w}))^2 \right] = \sum_{s \in S} \mu_\pi(s)[v_\pi(s) - \hat{v}(s; \mathbf{w})]^2 \tag{1.3}$$

The inclusion of $\mu_\pi(s)$ focuses the learning process on the parts of the state space that are most relevant to the agent's behavior.

## 1.4 The Learning Mechanism: Stochastic Gradient Descent (SGD)

With a differentiable objective function $J(\mathbf{w})$, we can employ **gradient descent** to find the parameters $\mathbf{w}$ that minimize the error. The gradient of the MSVE objective function with respect to the parameters $\mathbf{w}$ is:

$$\nabla_w J(\mathbf{w}) = \nabla_w \mathbb{E}_{s \sim \mu_\pi} \left[ (v_\pi(s) - \hat{v}(s; \mathbf{w}))^2 \right] \tag{1.4}$$

Applying the chain rule, this becomes:

$$\nabla_w J(\mathbf{w}) = -2\mathbb{E}_{s \sim \mu_\pi} \left[ (v_\pi(s) - \hat{v}(s; \mathbf{w}))\nabla_w \hat{v}(s; \mathbf{w}) \right] \tag{1.5}$$

The gradient descent update rule, scaled by a learning rate $\alpha$, is:

$$\Delta\mathbf{w} = -\frac{1}{2}\alpha\nabla_w J(\mathbf{w}) = \alpha\mathbb{E}_{s \sim \mu_\pi} \left[ (v_\pi(s) - \hat{v}(s; \mathbf{w}))\nabla_w \hat{v}(s; \mathbf{w}) \right] \tag{1.6}$$

In practice, we use **Stochastic Gradient Descent (SGD)**, which approximates the true gradient using a single sample at each time step $t$. The SGD update rule is:

$$\Delta\mathbf{w}_t = \alpha(v_\pi(S_t) - \hat{v}(S_t; \mathbf{w}))\nabla_w \hat{v}(S_t; \mathbf{w}) \tag{1.7}$$

This generalization is a double-edged sword. It allows RL to scale, but a poor update can corrupt the value function globally. This "crosstalk" between states is an intrinsic challenge of function approximation.

# Chapter 2

# Learning with Approximate Value Functions

## 2.1 The Target Problem: We Don't Know the True Value

The core challenge in applying SGD to RL is the absence of the ground-truth target $v_\pi(S_t)$. The solution is to substitute this unknown value with an *estimate*, or a *target*, computed from the agent's experience. The general update rule becomes:

$$\Delta\mathbf{w}_t = \alpha(\text{Target}_t - \hat{v}(S_t; \mathbf{w}))\nabla_w\hat{v}(S_t; \mathbf{w}) \tag{2.1}$$

The term $(\text{Target}_t - \hat{v}(S_t; \mathbf{w}))$ is the prediction error.

## 2.2 Prediction with Monte Carlo (MC) Targets

Monte Carlo (MC) methods use the actual, observed total discounted reward from a complete episode as the target. The **MC target** is the full return, $G_t$:

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots + \gamma^{T-t-1} R_T \tag{2.2}$$

The **MC update rule** is:

$$\Delta\mathbf{w}_t = \alpha(G_t - \hat{v}(S_t; \mathbf{w}))\nabla_w\hat{v}(S_t; \mathbf{w}) \tag{2.3}$$

**Properties of MC Learning:**

- **Unbiased Target:** $G_t$ is an unbiased sample of $v_\pi(S_t)$.

- **High Variance:** $G_t$ is a sum of many random variables, making it very noisy.

- **Episodic Requirement:** Updates can only be performed at the end of an episode.

## 2.3 Prediction with Temporal Difference (TD) Targets

Temporal Difference (TD) learning forms a target after only one step by **bootstrapping**: using the current value estimate to stand in for future returns. The **TD(0) target** is:

$$\text{TD Target} = R_{t+1} + \gamma\hat{v}(S_{t+1}; \mathbf{w}) \tag{2.4}$$

The **TD(0) update rule** is:

$$\Delta \mathbf{w}_t = \alpha(R_{t+1} + \gamma \hat{v}(S_{t+1}; \mathbf{w}) - \hat{v}(S_t; \mathbf{w}))\nabla_w \hat{v}(S_t; \mathbf{w}) \tag{2.5}$$

The term $\delta_t = R_{t+1} + \gamma \hat{v}(S_{t+1}; \mathbf{w}) - \hat{v}(S_t; \mathbf{w})$ is the **TD error**. **Properties of TD Learning:**

- **Biased Target:** The target depends on the current, potentially inaccurate, estimate $\hat{v}(S_{t+1}; \mathbf{w})$.

- **Lower Variance:** The target depends on only one reward and transition, making it less noisy.

- **Online and Incremental:** Updates can happen after every step.

Because the target itself depends on the parameters $\mathbf{w}$ and we ignore this dependency in the gradient calculation, these are called **semi-gradient methods**.

## 2.4 Control with Action-Value Function Approximation

To find the optimal policy (the control problem), we learn the action-value function $\hat{Q}(s, a; \mathbf{w})$. The policy is then derived by acting greedily: $a_t = \arg\max_a \hat{Q}(S_t, a; \mathbf{w})$.

- **"Action-in" Architecture:** Both state $s$ and action $a$ are inputs. This is necessary for **continuous action spaces**.

- **"Action-out" Architecture:** Only state $s$ is input; the network outputs a Q-value for each discrete action. This is used by DQN.

## 2.5 On-Policy vs. Off-Policy Updates

- **On-Policy (SARSA):** The agent learns the value of its current behavior policy. The target uses the action, $A_{t+1}$, that the agent *actually takes*.

$$\Delta \mathbf{w} = \alpha(R_{t+1} + \gamma \hat{Q}(S_{t+1}, A_{t+1}; \mathbf{w}) - \hat{Q}(S_t, A_t; \mathbf{w}))\nabla_w \hat{Q}(S_t, A_t; \mathbf{w}) \tag{2.6}$$

- **Off-Policy (Q-Learning):** The agent learns the value of the optimal (greedy) policy, regardless of the action it actually takes. The target uses the maximum possible Q-value.

$$\Delta \mathbf{w} = \alpha \left( R_{t+1} + \gamma \max_{a'} \hat{Q}(S_{t+1}, a'; \mathbf{w}) - \hat{Q}(S_t, A_t; \mathbf{w}) \right) \nabla_w \hat{Q}(S_t, A_t; \mathbf{w}) \tag{2.7}$$

# Chapter 3

# Deep Q-Networks (DQN)

## 3.1 The Instability of Deep RL: The "Deadly Triad"

Naively combining Q-learning with a deep neural network is notoriously unstable. This instability is often attributed to the **"deadly triad"**:

1. **Function Approximation:** Using a non-linear approximator like a neural network, where updates are not localized.

2. **Bootstrapping:** Updating estimates based on other estimates, which can propagate errors.

3. **Off-policy Learning:** A mismatch between the data collection policy and the target policy.

DQN introduced two key mechanisms to solve the primary challenges arising from this triad:

1. **Non-IID Data:** Solved by **Experience Replay**.

2. **Non-Stationary Targets:** Solved by a **Fixed Target Network**.

## 3.2 Solution 1: Experience Replay

The agent stores its experiences $e_t = (s_t, a_t, r_{t+1}, s_{t+1})$ in a large replay buffer. For training, it samples a random mini-batch of transitions from this buffer. **Benefits:**

- **Breaking Temporal Correlations:** Random sampling makes the training data closer to the IID (Independent and Identically Distributed) assumption of deep learning.

- **Increasing Sample Efficiency:** Each experience can be reused for multiple updates.

- **Smoothing Learning:** Updates are averaged over many past behaviors, preventing overfitting to recent experiences.

Experience replay transforms the RL problem into a sequence of standard supervised learning problems, making it compatible with deep learning techniques.

## 3.3   Solution 2: Fixed Q-Targets

The "moving target" problem arises because the target $y_i = r_i + \gamma \max_{a'} \hat{Q}(s_{i+1}, a'; \mathbf{w})$ depends on the same weights $\mathbf{w}$ that are being updated. DQN solves this by using two networks:

1. An **online network** with parameters $\mathbf{w}$, which is actively trained.

2. A **target network** with parameters $\mathbf{w}^-$, which is a periodic clone of the online network.

The target is calculated using the fixed target network:

$$y_i = r_i + \gamma \max_{a'} \hat{Q}(s_{i+1}, a'; \mathbf{w}^-) \tag{3.1}$$

The loss is computed between this stable target and the online network's prediction, $\hat{Q}(s_i, a_i; \mathbf{w})$. The target network weights $\mathbf{w}^-$ are updated by copying $\mathbf{w}$ every $C$ steps. This provides a stable, stationary target for the online network to learn from.

## 3.4 The Complete DQN Algorithm

---

**Algorithm 1** Deep Q-Network (DQN) with Experience Replay and Fixed Q-Targets

---

1: **Initialize:** Replay memory buffer $D$ with capacity $N$.
2: **Initialize:** Online Q-network $\hat{Q}$ with random weights $\mathbf{w}$.
3: **Initialize:** Target Q-network $\hat{Q}^-$ with weights $\mathbf{w}^- = \mathbf{w}$.
4: **Set hyperparameters:** learning rate $\alpha$, discount factor $\gamma$, exploration rate $\epsilon$, mini-batch size $B$, target network update frequency $C$.
5: **for** episode $= 1$ to M **do**
6:     Initialize environment and get initial state $s_1$.
7:     **for** t $= 1$ to T **do**
8:         **Action Selection:** With probability $\epsilon$, select a random action $a_t$.
9:         Otherwise, select $a_t = \arg\max_a \hat{Q}(s_t, a; \mathbf{w})$.
10:         **Environment Interaction:** Execute action $a_t$, observe reward $r_{t+1}$ and next state $s_{t+1}$.
11:         **Store Transition:** Store $(s_t, a_t, r_{t+1}, s_{t+1})$ in $D$.
12:         **Sample and Train:**
13:         Sample a random mini-batch of $B$ transitions $(s_j, a_j, r_{j+1}, s_{j+1})$ from $D$.
14:         For each transition $j$ in the mini-batch, compute the target $y_j$:
15:         **if** $s_{j+1}$ is terminal **then**
16:             $y_j = r_{j+1}$
17:         **else**
18:             $y_j = r_{j+1} + \gamma \max_{a'} \hat{Q}^-(s_{j+1}, a'; \mathbf{w}^-)$
19:         **end if**
20:         **Perform Gradient Descent:** Perform a gradient descent step on the online network parameters $\mathbf{w}$ to minimize the loss:
21:         $L(\mathbf{w}) = \frac{1}{B} \sum_{j=1}^{B} (y_j - \hat{Q}(s_j, a_j; \mathbf{w}))^2$
22:         $\Delta\mathbf{w} = \alpha \frac{1}{B} \sum_{j=1}^{B} (y_j - \hat{Q}(s_j, a_j; \mathbf{w})) \nabla_w \hat{Q}(s_j, a_j; \mathbf{w})$
23:         **Update Target Network:** Every $C$ steps, set $\mathbf{w}^- \leftarrow \mathbf{w}$.
24:     **end for**
25: **end for**

---

## 3.5 Empirical Validation: Ablation Study

An ablation study shows the impact of each component.

Table 3.1: Ablation Study of DQN Components on Atari Games (Scores)

| Game | Linear | Deep Net (Naive) | DQN w/ Targets Only | DQN w/ Replay Only |
|---|---|---|---|---|
| Breakout | 3 | 3 | 10 | 241 |
| Enduro | 62 | 29 | 141 | 831 |
| River Raid | 2345 | 1453 | 2868 | 4102 |
| Seaquest | 656 | 275 | 1003 | 823 |
| Space Invaders | 301 | 302 | 373 | 826 |

**Analysis of Results:**

- A naive deep network often performs worse than a linear model, highlighting the instability.

- Experience replay provides the most significant performance boost, suggesting that breaking data correlations is critical.

- The full DQN, combining both innovations, achieves the best performance, demonstrating a powerful synergistic effect.

# Part II

# Direct Policy Optimization

# Chapter 4

# The Policy Gradient Theorem

## 4.1 A New Paradigm: Directly Parameterizing the Policy

Policy-based methods directly parameterize the policy itself, $\pi_\theta(a|s)$, a function with parameters $\theta$ that maps states to a probability distribution over actions. **Advantages:**

- **Better Convergence Properties:** Smoother updates lead to more stable convergence.

- **Effectiveness in Continuous Action Spaces:** Naturally handles continuous actions by outputting parameters of a probability distribution (e.g., mean and variance of a Gaussian).

- **Ability to Learn Stochastic Policies:** Essential for partially observable environments or competitive games.

## 4.2 The Objective Function: Maximizing Expected Reward

The goal is to find parameters $\theta$ that maximize the expected total reward of a trajectory $\tau$. The objective function $J(\theta)$ is:

$$J(\theta) = \mathbb{E}_{\tau \sim p_\theta(\tau)}[r(\tau)] = \mathbb{E}_{\tau \sim p_\theta(\tau)}\left[\sum_{t=0}^{T-1} r(s_t, a_t)\right] \tag{4.1}$$

where $p_\theta(\tau)$ is the probability of a trajectory under policy $\pi_\theta$.

## 4.3 The Log-Derivative Trick

The main challenge is that $\theta$ influences the distribution $p_\theta(\tau)$ over which the expectation is taken.

$$\nabla_\theta J(\theta) = \nabla_\theta \int p_\theta(\tau) r(\tau) d\tau \tag{4.2}$$

We use the **log-derivative trick**, $\nabla_\theta f(x) = f(x)\nabla_\theta \log f(x)$, to rewrite the gradient as an expectation we can sample.

$$\nabla_\theta p_\theta(\tau) = p_\theta(\tau)\nabla_\theta \log p_\theta(\tau) \tag{4.3}$$

Substituting this back gives a tractable form:

$$\nabla_\theta J(\theta) = \int p_\theta(\tau)[\nabla_\theta \log p_\theta(\tau) r(\tau)]d\tau = \mathbb{E}_{\tau \sim p_\theta(\tau)}[(\nabla_\theta \log p_\theta(\tau))r(\tau)] \tag{4.4}$$

## 4.4   Derivation of the Policy Gradient Theorem

The log-probability of a trajectory is:

$$\log p_\theta(\tau) = \log p(s_0) + \sum_{t=0}^{T-1} \log \pi_\theta(a_t|s_t) + \sum_{t=0}^{T-1} \log p(s_{t+1}|s_t, a_t) \tag{4.5}$$

The gradient with respect to $\theta$ only affects the policy term:

$$\nabla_\theta \log p_\theta(\tau) = \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t|s_t) \tag{4.6}$$

This yields the final form of the **Policy Gradient Theorem**:

$$\nabla_\theta J(\theta) = \mathbb{E}_{\tau \sim p_\theta(\tau)}\left[\left(\sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t|s_t)\right)\left(\sum_{t=0}^{T-1} r(s_t, a_t)\right)\right] \tag{4.7}$$

**Interpretation:** The update increases the log-probability of actions taken in a trajectory, scaled by the total reward of that trajectory. Good trajectories are reinforced; bad ones are suppressed.

# Chapter 5

# REINFORCE: A Foundational Policy Gradient Algorithm

## 5.1 A Monte Carlo Approach

REINFORCE is a Monte Carlo algorithm that implements the Policy Gradient Theorem. It collects a full episode before performing an update. For a single episode (trajectory $\tau$), the gradient estimate is:

$$\nabla_\theta J(\theta) \approx \left( \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t|s_t) \right) \cdot r(\tau) \tag{5.1}$$

where $r(\tau)$ is the total reward for the episode.

## 5.2 The REINFORCE Algorithm

---

**Algorithm 2** REINFORCE (Monte Carlo Policy Gradient)

---

1: **Initialize:** Policy network $\pi_\theta$ with random parameters $\theta$.
2: **Set hyperparameters:** learning rate $\alpha$.
3: **for** episode i = 1 to M **do**
4:     **Generate an Episode:** Run policy $\pi_\theta$ to generate a trajectory $\tau_i = (s_0, a_0, r_1, \ldots, s_{T-1}, a_{T-1}, r_T)$.
5:     **Calculate Returns:** For each time step $t = 0, \ldots, T-1$, calculate the return-from-time-t:
6:       $G_t = \sum_{k=t+1}^{T} \gamma^{k-t-1} r_k$
7:       *(Note: Using the return-from-time-t, $G_t$, for each action $a_t$ is a crucial variance reduction technique.)*
8:     **Update Policy Parameters:** Loop through each time step $t$ of the episode and update $\theta$:
9:       $\theta \leftarrow \theta + \alpha \cdot G_t \cdot \nabla_\theta \log \pi_\theta(a_t|s_t)$
10:       *(In practice, updates are often batched over the episode or multiple episodes.)*
11: **end for**

---

## 5.3  Analysis of Strengths and Weaknesses

**Strengths:**

- **Simplicity and Generality:** Easy to implement and applies to continuous spaces.

- **Unbiased Gradient Estimate:** Guaranteed to move in the correct direction on average.

**Weaknesses:**

- **High Variance:** The Monte Carlo return $G_t$ is very noisy, leading to unstable training and slow convergence. This is the primary drawback.

- **Sample Inefficiency:** As an on-policy, Monte Carlo algorithm, it must discard all data after each update, requiring a constant stream of new trajectories.

# Part III

# A Comparative Synthesis of RL Paradigms

# Chapter 6

# Value-Based vs. Policy-Based Methods

## 6.1 Core Distinction: Learning "How Good" vs. "What to Do"

- **Value-Based (DQN):** Learns a value function $\hat{Q}(s, a; \mathbf{w})$ ("how good" is an action). The policy is *implicit* and derived via $\arg\max$.

- **Policy-Based (REINFORCE):** Learns a policy $\pi_\theta(a|s)$ ("what to do") directly. Value is used as a learning signal during training but is not stored.

## 6.2 Fundamental Trade-offs

- **Sample Efficiency:** Value-based methods are generally more sample-efficient due to off-policy learning and experience replay.

- **Stability and Convergence:** Policy-based methods are often more stable due to smoother updates.

- **Representational Complexity:** The optimal policy function is often simpler to represent with a neural network than the optimal value function.

## 6.3 Handling Action Spaces

- **Value-Based (DQN):** Best for **discrete action spaces** where the $\arg\max$ operation is tractable.

- **Policy-Based:** Naturally handles **both discrete and continuous action spaces**.

## 6.4 Nature of the Learned Policy

- **Value-Based (DQN):** Learns a **deterministic** policy. Stochasticity is added externally for exploration (e.g., $\epsilon$-greedy).

- **Policy-Based:** Naturally learns **stochastic** policies, which is crucial for partially observable or competitive environments.

# Chapter 7

# On-Policy vs. Off-Policy Learning

## 7.1 Defining the Concepts

- **Behavior Policy (b):** The policy used to generate experience.

- **Target Policy ($\pi$):** The policy being learned and improved.

**On-Policy Learning:** The policies are the same ($b = \pi$). The agent learns from its own actions. Examples: **SARSA**, **REINFORCE**. **Off-Policy Learning:** The policies can be different ($b \neq \pi$). The agent can learn about an optimal policy while behaving exploratorily. Example: **Q-Learning (DQN)**.

## 7.2 Algorithmic Implications

- **Data Utilization:** Off-policy methods are vastly more sample-efficient because they can reuse old data via **experience replay**. On-policy methods must discard data after each policy update.

- **Exploration:** Off-policy learning decouples exploration from control, allowing for more thorough exploration without compromising the learned optimal policy.

# Chapter 8

# Learning from Data: Reinforcement vs. Imitation

## 8.1 Behavioral Cloning: A Supervised Approach

**Behavioral Cloning (BC)** is a form of imitation learning. Given a dataset of expert state-action pairs $(s, a)$, it trains a policy $\pi_\theta(a|s)$ to mimic the expert. This is a standard **supervised learning** problem where states are inputs and actions are labels. The objective is to maximize the log-likelihood of the expert's actions:

$$\nabla_\theta J_{ML}(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \tag{8.1}$$

## 8.2 Contrasting Learning Signals

Comparing the policy gradient and behavioral cloning objectives reveals a deep connection:

- **Policy Gradient (REINFORCE):** $\nabla_\theta J_{PG}(\theta) \approx \mathbb{E}\left[\left(\sum_t \nabla_\theta \log \pi_\theta(a_t|s_t)\right) \cdot \left(\sum_t r_t\right)\right]$

- **Behavioral Cloning (ML):** $\nabla_\theta J_{ML}(\theta) \approx \mathbb{E}\left[\sum_t \nabla_\theta \log \pi_\theta(a_t|s_t)\right]$

**Insight:** Behavioral Cloning is mathematically equivalent to a policy gradient algorithm where the reward is always +1 for every action taken by the expert. It learns to mimic, not to optimize an external reward signal.

## 8.3 Applicability

- **Imitation Learning (BC)** is suitable when:

  - High-quality expert data is available.
  - Environmental interaction is expensive or dangerous.
  - A major challenge is **distribution shift**: small errors can lead the agent to unfamiliar states where its policy is poor.

- **Reinforcement Learning (RL)** is preferred when:

– An agent can interact with the environment safely and cheaply.

– A clear reward signal can be designed.

– The goal is to discover a policy superior to any human expert.

# Chapter 9

# Conclusion: Unifying Themes and Future Directions

This monograph has traversed the foundational landscape of modern reinforcement learning, revealing a field defined by a series of fundamental trade-offs. The choice between value-based and policy-based methods hinges on a trade-off between the **sample efficiency** of off-policy value learning and the **stability** of direct policy optimization. The success of landmark algorithms like DQN lies in their explicit mitigation of these issues through mechanisms like experience replay and fixed Q-targets.

The limitations of both pure paradigms have naturally led the field toward hybrid approaches that seek to combine their respective strengths. The most prominent of these are **Actor-Critic methods**. In an actor-critic architecture:

- The **Actor** is a policy network, $\pi_\theta(a|s)$, that learns "what to do," benefiting from the stability and continuous action space capabilities of policy-based methods.

- The **Critic** is a value network, $\hat{v}(s; \mathbf{w})$ or $\hat{Q}(s, a; \mathbf{w})$, that learns "how good" the actor's actions are. The critic's output is used to provide a low-variance learning signal to the actor, addressing the primary weakness of REINFORCE.

Actor-Critic methods represent a powerful synthesis, combining the direct policy optimization of policy gradients with the low-variance, bootstrapped learning signal from value-based methods. They form the backbone of many state-of-the-art algorithms today and represent a key step toward developing more robust, efficient, and general reinforcement learning agents.