



John Smith, Ph.D.

Department of Computer Science  
University Name

Email: [email@university.edu](mailto:email@university.edu)  
Website: [www.university.edu](http://www.university.edu)

July 13, 2025



John Smith, Ph.D.

Department of Computer Science  
University Name

Email: [email@university.edu](mailto:email@university.edu)  
Website: [www.university.edu](http://www.university.edu)

July 13, 2025

# Overview of Policy Gradient Methods

## What are Policy Gradient Methods?

Policy gradient methods are algorithms in reinforcement learning that optimize the policy directly by mapping states to actions and improving the policy based on gradients of expected rewards.

## Why are Policy Gradient Methods Important?

- **Direct Steering:** Suitable for continuous action spaces in complex environments like robotics and gaming.
- **High-Dimensional Action Spaces:** Effective where optimal actions are difficult to determine.
- **Flexibility:** Can be combined with value function approaches to enhance overall learning.

# How do Policy Gradient Methods Work?

## 1 Policy Representation:

- Denoted as  $\pi(a|s)$ : probability of taking action  $a$  in state  $s$ .
- Can be deterministic or stochastic.

## 2 Objective Function:

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_{t=0}^T r_t \right] \quad (1)$$

where  $\tau$  is a trajectory and  $r_t$  is the reward at time  $t$ .

## 3 Gradient Estimate:

$$\nabla J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_{t=0}^T \nabla \log \pi_{\theta}(a_t | s_t) R(\tau) \right] \quad (2)$$

$R(\tau)$  is the total return following actions taken by the policy.

# Key Points and Application

## Key Points

- **Versatility:** Adaptable for both discrete and continuous action spaces.
- **End-to-End Training:** Policies can be trained independently across episodes.
- **Drawbacks:** Variance in gradient estimates; often requires techniques like baselines or actor-critic methods for stabilization.

## Example Application

Consider a robot navigating a maze using policy gradient methods, exploring various actions and adjusting its policy based on rewards (e.g., reaching goals).

# Understanding Policy Gradients - Overview

- Definition of policy gradients
- Contrast with value-based methods
- Key concepts and motivations

# What are Policy Gradients?

- Policy Gradient methods optimize the policy directly.
- They parameterize the policy and improve it using gradient ascent based on rewards received.

## Key Concept: Policy

- **Policy ( $\pi$ ):** A function mapping states ( $s$ ) to probabilities of actions ( $a$ ).
- Example:  $\pi(a|s)$  is the probability of taking action  $a$  in state  $s$ .

# Motivation for Policy Gradients

- Suitable for large or continuous action spaces where value-based methods face challenges.
- Facilitate effective exploration and exploitation in complex environments.

## Contrasting with Value-Based Methods

- **Value-Based Methods:** Focus on estimating value functions (e.g., Q-values).
- Examples include Q-Learning and Deep Q-Networks (DQN).
- Optimization is done via maximizing the value function using techniques like Temporal Difference (TD) learning.



# Policy Gradient Methods Overview

- **Definition:** Parameterize the policy directly, optimizing it by maximizing expected returns using gradient ascent.
- **Examples:** REINFORCE algorithm, Proximal Policy Optimization (PPO).
- **Objective Function:**

$$J(\theta) = E_{\tau \sim \pi_{\theta}}[R(\tau)] \quad (3)$$

where  $R(\tau)$  is the total reward of trajectory  $\tau$ , and  $\theta$  are the policy parameters.

# Advantages and Challenges of Policy Gradients

## Advantages

- Suitable for continuous action spaces.
- Allow for stochastic policies, representing uncertainty in action selection.
- Directly optimize high-dimensional action spaces.

## Challenges

- High variance in gradient estimations, often needing variance reduction techniques.
- Lower sample efficiency compared to value-based methods.

# Policy Gradient Formula

$$\nabla J(\theta) = E_{s_t \sim \pi_\theta} [\nabla \log \pi_\theta(a_t | s_t) R_t] \quad (4)$$

- Where  $R_t$  is the total expected return from time  $t$  onward.
- $\nabla \log \pi_\theta(a_t | s_t)$  is the gradient of the log probability of the chosen action.

## Conclusion and Key Takeaway

- Policy gradient methods represent a powerful framework in reinforcement learning.
- They focus on optimizing policies directly, beneficial for complex and high-dimensional problems.
- Despite challenges like variance and sample inefficiency, mastering these methods is crucial for reinforcement learning.

### Key Takeaway

Policy gradients enable a direct approach to policy optimization, essential for tackling various challenges in reinforcement learning.

# Theoretical Foundations - Part 1

## Introduction to Policy Gradient Methods

- Policy gradient methods optimize policies directly in reinforcement learning.
- Focus on maximizing expected cumulative reward by adjusting policy parameters.
- Differ from value-based methods that estimate value functions.

# Theoretical Foundations - Part 2

## Mathematical Formulation

### ■ Expected Rewards:

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

- $J(\theta)$ : expected reward.
- $\tau$ : trajectory (sequence of states and actions).
- $R(\tau)$ : total reward over trajectory.

### ■ Policy Objective Function:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} [R(\tau) \nabla_{\theta} \log \pi_{\theta}(a|s)]$$

- How changes in  $\theta$  affect expected rewards.
- $\nabla_{\theta} \log \pi_{\theta}(a|s)$ : score function indicating parameter adjustment direction.

# Theoretical Foundations - Part 3

## Key Points

- Direct policy optimization vs value-based methods.
- Exploration vs exploitation in policy adjustment.
- Variance reduction techniques: reward normalization, baselines.

## Example: Visualizing the Objective Function

- Simple environment with actions A and B.
- If action A yields higher rewards than B, gradient indicates increasing action A's probability.

# REINFORCE Algorithm - Overview

## Overview

The REINFORCE algorithm is a foundational method in Reinforcement Learning (RL) representing a type of policy gradient method. This algorithm directly optimizes the policy by adjusting its parameters in the direction that increases expected rewards.



# REINFORCE Algorithm - Structure

- 1 **Initialization:** Start with a policy parameterized by  $\theta$ .
- 2 **Episode Generation:** Collect states, actions, and rewards:  $(s_t, a_t, r_t)$ .
- 3 **Return Calculation:** Compute the return  $G_t$  as:

$$G_t = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots \quad (5)$$

- 4 **Policy Update:** Update the policy parameters using:

$$\theta \leftarrow \theta + \alpha \nabla J(\theta) \quad (6)$$

where

$$\nabla J(\theta) = \frac{1}{T} \sum_{t=0}^T \nabla \log \pi(a_t | s_t; \theta) G_t \quad (7)$$

- 5 **Iteration:** Repeat until convergence.

# REINFORCE Algorithm - Key Points

- **Stochastic Policy:** Works with stochastic policies, outputting a distribution over actions.
- **Advantages:**
  - Simple to implement.
  - Directly optimizes expected rewards.
- **Limitations:**
  - High variance in returns, leading to unstable updates.
  - Sample inefficiency due to high episode requirements.
  - Lacks bootstrapping capabilities.

## REINFORCE Algorithm - Closing Thoughts

The REINFORCE algorithm serves as a stepping stone to more advanced policy gradient methods. It lays the groundwork for enhancements such as Actor-Critic approaches, which help mitigate the limitations observed in REINFORCE. Understanding this algorithm is crucial for further exploration in reinforcement learning techniques.

# Actor-Critic Methods - Overview

## Definition

Actor-Critic methods are a class of reinforcement learning algorithms that combine policy-based and value-based approaches.

- **Actor:** Selects actions based on the current policy.
- **Critic:** Evaluates the actions taken and provides feedback on their value.

# Actor-Critic Methods - Components

## Policy (Actor)

The policy is the strategy determining action selection, often represented as:

$$\pi(a|s; \theta) \tag{8}$$

where  $\theta$  are the parameters.

## Value Function (Critic)

The value function estimates the expected return from a state or state-action pair, typically represented by:

$$V(s; w) \tag{9}$$

where  $w$  denote the critic's parameters.

# Actor-Critic Methods - Learning Process

## 1 Interaction with Environment:

- The agent takes actions based on the actor and receives rewards.

## 2 Critic Feedback:

- Evaluates action using the advantage function:

$$A_t = r_t + \gamma V(s_{t+1}; w) - V(s_t; w) \quad (10)$$

## 3 Policy Update:

$$\theta \leftarrow \theta + \alpha \cdot A_t \nabla_{\theta} \log \pi(a_t | s_t; \theta) \quad (11)$$

## 4 Value Function Update:

$$L(w) = \frac{1}{2} (v_t - V(s_t; w))^2 \quad (12)$$

# Variations of Policy Gradient Methods

## Overview

Policy Gradient Methods are a class of reinforcement learning techniques that directly optimize the policy used to select actions. This presentation focuses on two prominent variations:

- **Proximal Policy Optimization (PPO)**
- **Trust Region Policy Optimization (TRPO)**

Both address key challenges in policy optimization.

# Proximal Policy Optimization (PPO)

## Concept

PPO balances performance and training stability using a clipped objective function to restrict policy updates.

## Clipped Objective Function

$$L^{CLIP}(\theta) = \mathbb{E}_t \left[ \min \left( r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right] \quad (13)$$

where:

- $r_t(\theta) = \frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$
- $\hat{A}_t$  is the estimated advantage at time  $t$ .
- $\epsilon$  controls the range of policy updates.

## Key Features



# Trust Region Policy Optimization (TRPO)

## Concept

TRPO uses a theoretical framework ensuring policy updates do not move too far from the current policy by constraining the KL divergence.

## Objective and Constraint

$$\text{maximize } \mathbb{E}_t [\hat{A}_t] \quad \text{subject to } \mathbb{E}_t [D_{KL}(\pi_{\theta_{old}} || \pi_{\theta})] \leq \delta \quad (14)$$

where  $D_{KL}$  measures divergence and  $\delta$  is a threshold.

## Key Features

- Conservative policy changes due to the "trust region."
- More computationally intensive; requires solving constrained optimization.
- Mitigates policy collapse and instability.

# Applications of Policy Gradients - Introduction

- Policy Gradient Methods are reinforcement learning algorithms that optimize policies by gradient ascent on expected rewards.
- Effective in complex environments where traditional value-based methods struggle.
- Capable of handling:
  - High-dimensional action spaces
  - Continuous action spaces
  - Partially observable states

# Applications of Policy Gradients - Key Applications

## 1 Robotics

- Example: Training robotic arms to grasp objects.
- Policy gradients allow learning intricate tasks through trial and error.

## 2 Game Playing

- Example: Playing games like Atari and Go.
- State-of-the-art performance achieved via self-play and exploration through algorithms like PPO and TRPO.

## 3 Autonomous Driving

- Example: Navigation in complex environments.
- Learn driving policies by interacting with realistic simulations.

## 4 Finance

- Example: Portfolio management.
- Reinforcement learning agents learn trading strategies based on market simulations.

## 5 Natural Language Processing

- Example: Dialogue systems and machine translation.
- Policy gradients help maximize user engagement and contextual relevance.

# Applications of Policy Gradients - Importance and Challenges

## Importance and Advantages

- **Flexibility:** Handles stochastic policies for exploration in uncertain environments.
- **Direct Policy Optimization:** Avoids complications from value function estimation.
- **Scalability:** Suitable for high-dimensional spaces without easy discretization.

## Challenges and Considerations

- **High Variance:** Estimates of gradients can lead to unstable training; techniques like Baseline subtraction help.
- **Sample Inefficiency:** Requires significant training data which can be costly.

# Implementation of Policy Gradient Algorithms

## Overview

This presentation covers the implementation of a policy gradient method in reinforcement learning. Key topics include:

- Understanding Policy Gradient Methods
- Key Concepts
- Implementation Steps using Python (PyTorch)
- Summary and Key Points

# Understanding Policy Gradient Methods

- Policy Gradient methods optimize the policy directly.
- They adjust parameters based on expected returns.
- Unlike value-based methods, these methods choose actions based on a parameterized probability distribution.

# Key Concepts

- 1 Policy Representation:** A policy  $\pi$  is represented as a neural network. Input: state, Output: action probabilities.
- 2 Objective Function:** The aim is to maximize the expected return

$$J(\theta) = E[R] \quad (15)$$

where  $R$  is the return from a state following the policy.

# Implementation Steps in Python (Part 1)

## 1 Install Required Libraries:

```
pip install torch gym
```

## 2 Environment Setup:

```
import gym
env = gym.make('CartPole-v1') # Example environment
```

## 3 Define the Policy Network:

```
import torch.nn as nn
import torch.optim as optim

class PolicyNetwork(nn.Module):
    def __init__(self, input_dim, output_dim):
        super(PolicyNetwork, self).__init__()
        self.fc = nn.Linear(input_dim, 128)
```



## Implementation Steps in Python (Part 2)

### 4 Define the Training Function:

```
def train(env, policy_net, optimizer, num_episodes=1000,
        for episode in range(num_episodes):
            state = env.reset()
            rewards = []
            log_probs = []

            for step in range(max_steps):
                state_tensor = torch.FloatTensor(state)
                action_probs = policy_net(state_tensor)
                action = torch.multinomial(action_probs, 1).i

            # Sampling action

            next_state, reward, done, _ = env.step(action
```

## Implementation Steps in Python (Part 3)

### 5 Compute Returns:

```
def compute_returns(rewards, discount_factor=0.99):  
    returns = []  
    R = 0  
    for r in reversed(rewards):  
        R = r + discount_factor * R  
        returns.insert(0, R)  
    return torch.FloatTensor(returns)
```

### 6 Initialization and Execution:

```
input_dim = env.observation_space.shape[0]  
output_dim = env.action_space.n  
policy_net = PolicyNetwork(input_dim, output_dim)  
optimizer = optim.Adam(policy_net.parameters(), lr=0.01)
```

## Key Points and Summary

- Policy Gradient Basics: Focus on adjusting policy parameters using gradients from sampled returns.
- Neural Network Structure: A neural network enables flexible policy approximation.
- Exploration vs. Exploitation: Important to balance exploring diverse actions and exploiting known benefits.

### Summary

Implementing policy gradient methods involves:

- defining a policy,
- collecting trajectories,
- updating the policy based on the returns received.

This serves as a foundation for advanced reinforcement learning techniques.

# Performance Evaluation - Introduction

Performance evaluation of policy gradient methods is crucial to understand their effectiveness in solving reinforcement learning (RL) problems. This section discusses:

- Key metrics for assessment
- Sampling efficiency
- Convergence properties

# Performance Evaluation - Key Metrics

## 1 Average Return:

- Measures cumulative reward over episodes.
- Formula:

$$R = \frac{1}{N} \sum_{i=1}^N R_i \quad (16)$$

(where  $R_i$  is the return of episode  $i$  and  $N$  is total episodes)

## 2 Sample Efficiency:

- Evaluates number of samples needed for performance level.
- Higher efficiency indicates quicker learning.

## 3 Convergence Properties:

- Speed and reliability of reaching optimal solutions.
  - **Rate of Convergence:** How fast the algorithm approaches optimal policy.
  - **Stability:** Resilience against variations.

# Performance Evaluation - Techniques

Techniques for evaluating policy gradient methods:

## 1 Multiple Experiments:

- Conduct runs with different random seeds and average results.

## 2 Learning Curves:

- Visualize average return vs. number of episodes.
- *Python Code Example:*

```
import matplotlib.pyplot as plt

episodes = [i for i in range(len(average_returns))]
plt.plot(episodes, average_returns)
plt.title('Learning Curve')
plt.xlabel('Episodes')
plt.ylabel('Average Return')
plt.show()
```

## 3 Hyperparameter Sensitivity Analysis:

# Challenges and Future Directions in Policy Gradient Methods

## Overview

Policy Gradient Methods (PGMs) have shown great promise in reinforcement learning (RL) but face several challenges. This presentation discusses the hurdles in deploying PGMs effectively and outlines future research directions to enhance their performance.

# Current Challenges

## 1 High Variance in Gradient Estimates

- Explanation: High variance can lead to unstable learning and convergence issues.
- Example: The REINFORCE algorithm offers unbiased estimates with high variance, causing slow convergence.
- Mitigation: Variance reduction techniques (e.g., using baselines).

## 2 Sample Inefficiency

- Explanation: Requires a large number of samples, which can be computationally expensive.
- Example: Limited data environments make sufficient sample acquisition impractical.
- Future Direction: Explore methods to enhance sample efficiency using advanced function approximation.



## Current Challenges (continued)

### 3 Suboptimal Local Optima

- Explanation: PGMs may converge to suboptimal policies due to local minima.
- Example: Complex environments may trap the policy in a less optimal strategy.
- Future Direction: Investigate methods like simulated annealing to escape local minima.

### 4 Combining Exploration and Exploitation

- Explanation: Balancing new action exploration and leveraging known actions is challenging.
- Example: Over-exploitation can lead to missing better strategies.
- Future Direction: Integrate exploration strategies (e.g.,  $\epsilon$ -greedy methods) to improve learning.

### 5 Complexity of Policy Representation

- Explanation: Designing efficient architectures that generalize across states is difficult.
- Example: Deep reinforcement learning struggles with creating efficient networks while maintaining performance.
- Future Direction: Research on expressive policy representations, like hierarchical reinforcement learning.

# Future Research Directions

## 1 Hybrid Approaches

- Integrating PGMs with value-based methods to create robust frameworks.

## 2 Improved Algorithms

- Developing adaptive algorithms that adjust based on training performance feedback.

## 3 Addressing Computational Demand

- Optimizing performance on resource-constrained devices through lightweight models.

## 4 Robustness to Environment Changes

- Exploring transfer learning and meta-RL for quick adaptability to new tasks.

## Key Takeaways

- PGMs present unique challenges, including variance in gradient estimates and sample inefficiency.
- Future research should focus on hybrid methods, improved algorithms, and enhancing robustness to realize the full potential of PGMs in practical applications.

### Note

Understanding these challenges and future directions is vital for advancing the applications of policy gradient methods in various domains, from robotics to game AI.

# Ethical Considerations in RL - Introduction

- **Ethics in AI:** The integration of ethical considerations into AI technologies is crucial for positive societal impact.
- **Importance of Responsible AI:** Understanding the consequences of decisions made by RL systems is vital to prevent harm and promote fairness.

# Ethical Considerations in RL - Key Ethical Concerns

## 1 Bias and Fairness:

- RL algorithms may learn biased behaviors from training data or reward structures.
- *Example:* An RL agent in a game may favor certain strategies, leading to unfair experiences.

## 2 Transparency and Explainability:

- Complex models from policy gradient methods obscure understanding, affecting trust and accountability.
- *Example:* In accidents caused by self-driving cars, understanding decision-making is crucial.

## 3 Safety and Control:

- Ensuring safe operations in unpredictable environments is essential.
- *Example:* A robot arm optimizing for speed may risk human safety.

## 4 Long-term Consequences:

- RL agents should consider long-term impacts and potential adverse outcomes.
- *Example:* Maximizing short-term engagement in apps can lead to addiction.

# Ethical Considerations in RL - Best Practices

- **Define Clear Ethical Guidelines:** Establish a framework for ethical decision-making, addressing bias and transparency.
- **Continuous Monitoring and Evaluation:** Ongoing assessments to identify and correct ethical issues.
- **Stakeholder Involvement:** Engage diverse stakeholders for a comprehensive understanding of impacts.
- **Adopt Explainable AI (XAI) Techniques:** Techniques that improve explainability enhance user trust and accountability.

# Ethical Considerations in RL - Conclusion and Key Takeaways

- Ensuring ethical standards enhances trustworthiness and aligns AI with societal values.
- **Key Takeaways:**
  - Ethics is integral to RL framework design.
  - Transparency, fairness, and safety are critical ethical pillars.
  - Engaging with stakeholders enhances ethical outcomes.

## Conclusion - Recap of Policy Gradient Methods

- Policy Gradient methods are crucial in reinforcement learning (RL) as they directly optimize the policy.
- They differ from value-based methods by modeling the probability distribution over actions given states.



## Conclusion - Key Advantages of Policy Gradient Methods

- 1 **Continuous Action Spaces:** Effective in environments like robotic control.
- 2 **High-Dimensional Action Spaces:** Efficient in scenarios such as game playing and autonomous driving.
- 3 **Stochastic Policies:** Able to explore and exploit simultaneously, enhancing performance in uncertain environments.

## Conclusion - Example and Future Impact

### Example: REINFORCE Algorithm

The REINFORCE algorithm updates policy parameters  $\theta$  using:

$$\theta \leftarrow \theta + \alpha \cdot \nabla J(\theta)$$

where

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

### Future Advancements:

- **Healthcare:** Personalized treatment strategies.
- **Gaming:** Intelligent non-playable characters (NPCs).
- **Autonomous Systems:** Enhanced performance in dynamic environments.