# Week 6: Policy Gradient Methods

Your Name

Your Institution

June 30, 2025

# Introduction to Policy Gradient Methods

## Overview of Policy Gradients in Reinforcement Learning

Reinforcement Learning (RL) is a branch of machine learning where an agent learns to make decisions by interacting with an environment. The goal is to maximize cumulative rewards through a trial-and-error approach.

- **Agent**: The decision-maker that interacts with the environment.
- **Environment**: The external system with which the agent interacts.
- **State**: A snapshot of the environment at a particular time.
- **Action**: Choices made by the agent that affect the state.
- **Reward**: A signal received from the environment indicating the value of the action taken.

## Policy Gradient Methods

A type of RL technique that optimizes the policy directly, contrasting with value-based methods that estimate value functions.

# REINFORCE Algorithm

## Overview

REINFORCE is one of the most well-known policy gradient algorithms.

- **Stochastic Policy**: The policy is represented as a probability distribution over actions given a state.
- **Episode-Based**: The agent interacts with the environment until an episode ends and receives a total reward.
- **Gradient Estimation**: The update rule uses the observed actions and their rewards to improve the policy.

$$\theta' = \theta + \alpha \cdot \nabla J(\theta) \tag{1}$$

- $\theta$: parameters of the policy.
- $\alpha$: learning rate.
- $\nabla J(\theta)$: estimated gradient of the return.

## Return Calculation

# Example Scenario

## Grid Navigation Game

Consider a simple game where an agent navigates a grid environment to reach a goal:

- The agent can move up, down, left, or right.
- Each action may lead to positive or negative rewards based on proximity to the goal.

## Learning Process

After numerous episodes, the agent collects rewards (e.g., +10 for reaching the goal) and updates its policy to increase the likelihood of successful actions based on those experiences, utilizing REINFORCE.

# Key Points and Conclusion

## Key Points

- Reinforcement Learning focuses on learning from interactions with the environment.
- Policy Gradient methods, like REINFORCE, optimize policies directly.
- Understanding the REINFORCE algorithm is fundamental for grasping how agents learn effective behaviors.

## Conclusion

Policy Gradient Methods, and specifically the REINFORCE algorithm, are powerful tools in reinforcement learning, allowing agents to learn adaptive policies through exploration and experience.

# Fundamental Concepts in Reinforcement Learning - Overview

## Key Concepts

This slide covers essential concepts relevant to reinforcement learning and policy gradient methods:

- Agent
- Environment
- State
- Action
- Reward
- Value Function

# Fundamental Concepts in Reinforcement Learning - Agent and Environment

1. **Agent**
   - The learner or decision-maker in an RL scenario.
   - Interacts with the environment to achieve specific goals.
   - **Example:** A robot navigating through a maze.
2. **Environment**
   - Encompasses everything the agent interacts with.
   - Defines the context where the agent operates.
   - **Example:** The maze itself with walls and exit points.

# Fundamental Concepts in Reinforcement Learning - States, Actions, and Rewards

3. **State**
   - A specific configuration of the environment at a given time.
   - Represents the current situation observed by the agent.
   - **Example:** Agent's current position in the maze (e.g., coordinates (5, 3)).

4. **Action**
   - A decision made by the agent that affects the environment's state.
   - The set of possible actions defines the agent's behavior.
   - **Example:** Moving left, right, up, or down in the maze.

5. **Reward**
   - A scalar value received after taking an action in a specific state.
   - Serves as feedback shaping the agent's learning.
   - **Example:** +10 for reaching the exit, -1 for hitting a wall.

6. **Value Function**
   - Estimates how good a specific state (or state-action pair) is in terms of expected future rewards.
   - **Example:** The total expected rewards if starting in a certain state.

## Integration with Policy Gradient Methods

These concepts are foundational to policy gradient methods:

- The agent learns to optimize its policy based on rewards.
- Policy adjustment is distinct from value-based methods.

## Policy Gradient Estimation

$$\nabla J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[ \nabla \log \pi_\theta(a_t | s_t) R(\tau) \right] \tag{3}$$

# Fundamental Concepts in Reinforcement Learning - Conclusion

Understanding these fundamental concepts is essential for grasping how policy gradient methods function in reinforcement learning. The interactions between agents, environments, states, actions, and rewards lay the groundwork for exploring advanced topics in this field.

# Understanding Policy Gradients - Introduction

## What are Policy Gradient Methods?

Policy Gradient methods are a class of algorithms in Reinforcement Learning (RL) that optimize the policy directly to select actions based on the current state. They parameterize the policy function $\pi(a|s; \theta)$, where $\theta$ represents the parameters, $s$ is the state, and $a$ is the action.

# Understanding Policy Gradients - Key Concepts

- **Policy vs. Value:**
  - **Policy** ($\pi$): A mapping from states to action probabilities, defining the agent's behavior.
  - **Value Function**: Estimates the expected return of taking certain actions in a given state.

- **Objective Function:** The goal is to maximize the expected return:

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau)] \qquad (4)$$

where $R(\tau)$ is the total return from trajectory $\tau$.

- **Direct Action Selection:** - Optimizes action probabilities, allowing for stochastic policies that can better handle uncertainty.
- **Rich Action Spaces:** - Effective in environments with high-dimensional, continuous action spaces (e.g., robotics, video games).
- **Handling Large State Spaces:** - Works well in complex environments where the action-value function is difficult to define.
- **No Maximum Action Issues:** - Avoids the problem of selecting a maximum action as in Q-learning, improving exploration capabilities.

- **Complex Decision Making:** - Problems with long-term dependencies and delayed rewards, such as in partially observable environments.
- **When Actions are Stochastic:** - Adaptively adjust action probabilities in non-deterministic scenarios, like recommendation systems or games.

# Understanding Policy Gradients - Example

## Illustrative Example

Consider a robot navigating a maze:

- The robot can choose from multiple paths (actions) at each junction (state).
- A value-based method may lead to a single perceived optimal path, potentially missing better long-term rewards.
- A policy gradient approach allows it to explore different paths, adjusting actions based on success and failure.

# Understanding Policy Gradients - Conclusion

## Summary

Policy gradient methods are powerful in RL, especially for continuous action spaces and complex decision-making tasks. Understanding their mechanics provides a foundation for exploring specific algorithms, such as REINFORCE.

## Key Points to Remember

- They optimize actions directly and are advantageous in continuous domains.
- Allow exploration of stochastic policies, crucial in uncertainty and multi-path scenarios.
- Notable effectiveness in real-world applications like robotics and adaptive game AI.

# The REINFORCE Algorithm - Overview

- The REINFORCE algorithm is a fundamental policy gradient method.
- It helps agents learn optimal policies using episodic returns.
- Adjusts action selection probabilities based on rewards received.

1. **Policy**:
   - Denoted as $\pi(a|s)$, representing a probability distribution over actions $a$ given a state $s$.
   - The agent follows this policy to make decisions in the environment.

2. **Return**:
   - Calculated as:
   $$G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \ldots \tag{5}$$
   - Where $\gamma$ is the discount factor ($0 < \gamma < 1$).

- The objective is to optimize the expected return:

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [G_t] \tag{6}$$

- Update rule for policy parameters $\theta$:

$$\theta \leftarrow \theta + \alpha \nabla J(\theta) \tag{7}$$

- Gradient of the objective using the log-likelihood trick:

$$\nabla J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [G_t \nabla \log \pi_\theta(a_t | s_t)] \tag{8}$$

1. **Sample Trajectories**:
   - Generate episodes using the current policy $\pi_\theta$.
2. **Compute Returns**:
   - Compute return $G_t$ for each time step in the episode.
3. **Policy Update**:
   - Update using:

$$\Delta\theta = \frac{1}{T}\sum_{t=0}^{T} G_t \nabla \log \pi_\theta(a_t|s_t) \tag{9}$$

- Consider an agent moving in a simplified environment (left or right).
- If the agent moves right and collects a reward of 1, future expected rewards influence the return calculation.
- After an episode, the agent adjusts its policy to favor actions that led to greater returns.

# The REINFORCE Algorithm - Key Points

- Utilizes Monte Carlo sampling to estimate returns, suitable for episodic tasks.
- High variance in updates can be mitigated with:
  - Reward scaling
  - Using a baseline

- The REINFORCE algorithm is foundational in reinforcement learning.
- It optimizes policies directly, paving the way for advanced methods.
- Understanding its mechanics is crucial for implementing and innovating policy gradient techniques.

# Implementing REINFORCE - Overview

## Objective

This slide provides a comprehensive, step-by-step guide to implementing the REINFORCE algorithm using Python, particularly with TensorFlow or PyTorch.

## Overview of the REINFORCE Algorithm

REINFORCE is a Monte Carlo policy gradient method that optimizes a parameterized policy using the policy gradient theorem. It directly updates the policy parameters based on the generated episodes.

1. **Setup Environment:**
   - Install necessary libraries:
     ```
     pip install gym torch numpy
     ```
   - Import required modules:
     ```
     import numpy as np
     import torch
     import torch.nn as nn
     import torch.optim as optim
     import gym
     ```

2. **Define the Policy Network:**
   ```
   class PolicyNetwork(nn.Module):
       def __init__(self, input_size, output_size):
           super(PolicyNetwork, self).__init__()
           self.fc = nn.Linear(input_size, 128)   #
               Hidden layer
           self.fc_out = nn.Linear(128, output_size)
               # Output layer
   ```

③ **Initialize the Environment and Agent:**

```
env = gym.make('CartPole-v1')
policy_net = PolicyNetwork(input_size=env.
    observation_space.shape[0], output_size=env.
    action_space.n)
optimizer = optim.Adam(policy_net.parameters(), lr
    =0.01)
```

④ **Run Episodes and Collect Trajectories:**

```
def run_episode(env, policy_net):
    state = env.reset()
    log_probs = []
    rewards = []
    done = False
    while not done:
        state_tensor = torch.FloatTensor(state).
            unsqueeze(0)
        action_probs = policy_net(state_tensor)
```

# Analyzing Performance Metrics

- Discussion on performance evaluation metrics for policy gradient methods.
- Key topics include:
  - Cumulative reward
  - Convergence rates
  - Factors impacting performance

# Understanding Performance Metrics in Policy Gradient Methods

## 1. Cumulative Reward

- **Definition**: Total reward over episodes, a holistic measure of agent's performance.
- **Formula**:

$$C = \sum_{t=0}^{T} R_t$$

- **Example**: For rewards of 5, -1, 3, cumulative reward: $5 + (-1) + 3 = 7$.

## 2. Convergence Rates

- **Definition**: Speed at which learning algorithm approaches a stable policy.
- **Monitoring**: Plot average cumulative reward over episodes for visual indication of convergence.
- **Example**: Trend of averages over 100 episodes can depict effective convergence.

## 3. Factors Impacting Performance

- **Hyperparameters**: Learning rate, discount factor, batch size, architecture.
- **Exploration vs. Exploitation**: Use strategies like $\epsilon$-greedy or softmax action selection.
- **Variance Reduction Techniques**: Normalize rewards or use Generalized Advantage Estimation (GAE).

# Key Points and Example Code

## Key Points to Emphasize

- Performance metrics provide quantifiable measures of effectiveness for policy gradients.
- Cumulative reward serves as a foundational metric for tracking progress.
- Convergence rates give insights into how efficiently the algorithm learns.
- Hyperparameters and exploration/exploitation strategies directly impact results.

## Example Code Snippet for Cumulative Reward Calculation

```
# Assuming rewards is a list of rewards obtained in an
    episode
cumulative_reward = sum(rewards)
print("Cumulative Reward:", cumulative_reward)
```

# Comparing Policy Gradient Methods

## Overview

Policy Gradient Methods are a class of reinforcement learning algorithms that optimize the policy directly utilizing gradient ascent techniques. Various approaches have unique strengths and weaknesses.

1. **Basic Policy Gradient (REINFORCE)**
   - **Description**: Uses the log probability of actions weighted by returns.
   - **Strengths**:
     - Simple to implement.
     - Effective in environments with discrete action spaces.
   - **Weaknesses**:
     - High variance in updates, leading to slow convergence.
   - **Formula**:
     $$\nabla J(\theta) = \mathbb{E}[\nabla \log \pi_\theta(a_t|s_t) R_t] \qquad (10)$$

2. **Actor-Critic Methods**
   - **Description**: Combines a policy function (actor) and a value function (critic).
   - **Strengths**:
     - Faster convergence due to variance reduction from the critic.
     - More sample-efficient.
   - **Weaknesses**:
     - Requires careful tuning of two components; may cause instability.
   - **Illustration**: The actor computes action probabilities while the critic evaluates actions by estimating the value function.

3. **Trust Region Policy Optimization (TRPO)**
   - **Description**: Optimizes the policy within a trust region for stability.
   - **Strengths**:
     - Guarantees monotonic improvement.
     - More stable training with larger updates.
   - **Weaknesses**:
     - Computationally expensive due to constrained optimization needs.
   - **Key Concept**: Kullback-Leibler (KL) divergence limits policy updates.
4. **Proximal Policy Optimization (PPO)**
   - **Description**: An improvement over TRPO with a clipped objective function.
   - **Strengths**:
     - Balances exploration and exploitation effectively.
     - Easier implementation than TRPO while maintaining performance.
   - **Weaknesses**:
     - May be less stable than TRPO depending on hyperparameters.
   - **Formula**:

$$L(\theta) = \mathbb{E}\left[\min\left(\frac{\pi_\theta(a|s)}{\pi_{\theta_{old}}(a|s)}A_t, g(\epsilon, A_t)\right)\right] \tag{11}$$

# Summary of Strengths and Weaknesses

| Method | Strengths | Weaknesses |
|---|---|---|
| REINFORCE | Simple, suited for discrete actions | High variance, slow |
| Actor-Critic | Reduces variance, sample-efficient | Instability if actor/ |
| TRPO | Monotonic improvement, stable | Computationally he |
| PPO | Simpler than TRPO, good performance | May require carefu |

## Key Takeaways

- Policy Gradient Methods are vital in reinforcement learning to maximize expected rewards. - Each method's optimization, variance control, and stability dictate its effectiveness. - Understanding trade-offs allows practitioners to select the best approach for their applications.

# Conclusion

As we've explored, while the foundational principle of policy gradients remains consistent, the choice of method significantly influences performance and stability in complex environments. Next, we will look at real-world case studies showcasing these methods in action.

# Case Studies of Policy Gradient Applications

## Introduction to Policy Gradient Methods

Policy Gradient Methods are a class of reinforcement learning algorithms that optimize the policy directly. Unlike value-based methods that estimate the value function, policy gradients adjust the parameters of the policy by maximizing the expected reward through deferred feedback.

# Real-World Applications of Policy Gradient Methods

1. **Robotics**
   - **Example**: Robotic Hand Manipulation
   - **Context**: Training robotic arms for complex tasks.
   - **Key Insight**: Handles continuous action space, enabling smooth transitions in tasks.

2. **Finance**
   - **Example**: Algorithmic Trading
   - **Context**: Automated trading strategies in dynamic markets.
   - **Key Insight**: Adapts to market noise through probabilistic strategy optimization.

3. **Healthcare**
   - **Example**: Personalized Treatment Recommendations
   - **Context**: Suggesting tailored treatment plans using patient health records.
   - **Key Insight**: Continuously learns from patient responses to refine suggestions.

# Key Insights and Conclusions

## Key Points to Emphasize

- **Versatility**: Handles both discrete and continuous action spaces.
- **Real-time Learning**: Adapts based on direct interaction with environments.
- **Exploration vs. Exploitation**: Balances these components using stochastic policy representations.

## Conclusion

Policy Gradient Methods show remarkable promise across various fields. They foster adaptability and optimize actions based on learned experiences, shaping the future of automation and decision-making.

# Formulas in Policy Gradient Methods

## Policy Objective Function

$$J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum_{t=0}^{T} r_t \right] \tag{12}$$

Where $\tau$ is the trajectory of states and rewards, and $r_t$ is the reward at time $t$.

## Gradient Ascent Update

$$\theta_{new} = \theta_{old} + \alpha \nabla J(\theta) \tag{13}$$

Here, $\alpha$ is the learning rate.

## Understanding Ethical Implications

As reinforcement learning (RL) techniques, especially policy gradient methods, integrate further into various domains, it is crucial to examine the ethical implications, with bias and fairness being significant concerns.

- Bias in Policy Gradients:
  - **Definition:** Bias occurs when the learned policy reflects prejudiced assumptions from training data.
  - **Source of Bias:** Historical data can encode systemic inequalities (e.g., gender, racial bias).
- Fairness in Decision-Making:
  - **Definition:** Impartial treatment of individuals regardless of their characteristics.
  - **Importance:** Bias in policies can lead to unfair treatment, reinforcing existing inequalities, particularly in hiring or law enforcement.

# Ethical Considerations in Policy Gradients - Influencing Factors

1. **Training Data Quality:**
   - Example: An RL system trained on biased hiring data may favor certain demographic groups.
2. **Reward Structure:**
   - Example: A reward function that maximizes clicks may disadvantage specific user groups without inclusivity metrics.
3. **Exploration Strategies:**
   - Exploration versus exploitation can unintentionally favor certain outcomes, leading to unequal treatment.

## Introduction to Emerging Trends

Policy gradient methods are at the forefront of reinforcement learning (RL) advancements. This slide explores key directions in which policy gradient methods are heading, emphasizing emerging trends and research opportunities.

# Future Directions in Policy Gradients - Sample Efficiency Improvement

1. **Sample Efficiency Improvement**
   - **Concept**: Enhancing how effectively algorithms learn from limited data.
   - **Example**: Techniques like *experience replay* that allow models to revisit important past experiences more often.
   - **Key Point**: Aim for methods that reduce the number of required samples for training, making RL feasible in high-cost data scenarios.

2. **Hierarchical Reinforcement Learning**
   - **Concept**: Structuring the learning problem into levels of abstraction for faster, coherent learning.
   - **Example**: Using sub-policies for multiple tasks, allowing agents to achieve high-level goals while managing sub-goals.
   - **Key Point**: This approach simplifies complex tasks, enhancing policy training by mimicking human-like decision-making.

# Future Directions in Policy Gradients - Integration with Deep Learning

3. **Integration with Deep Learning**
   - **Concept**: Combining policy gradient methods with deep learning to handle high-dimensional state spaces.
   - **Example**: Using Convolutional Neural Networks (CNNs) for visual data processing in robotic control.
   - **Key Point**: Deep learning enhances policy gradient methods' performance in complex feature extraction environments.

# Future Directions in Policy Gradients - Continual Learning and Explainability

4. **Continual and Lifelong Learning**
   - **Concept**: Allowing agents to learn and adapt continuously from new experiences.
   - **Example**: An agent retains knowledge when trained on different games, utilizing transfer learning.
   - **Key Point**: Focuses on reducing catastrophic forgetting and enhancing robustness in dynamic environments.

5. **Explainability and Interpretability in RL**
   - **Concept**: Developing methods to understand agent decisions in reinforcement learning.
   - **Example**: Visualizing policy behavior to help users understand decision pathways.
   - **Key Point**: Addressing the "black box" nature of deep RL is crucial for trust and ethical considerations.

# Future Directions in Policy Gradients - Multi-Agent Learning and Conclusion

6. **Multi-Agent Reinforcement Learning (MARL)**
   - **Concept**: Learning strategies when multiple agents operate in the same environment.
   - **Example**: Training autonomous vehicles to navigate traffic considering other vehicles' behavior.
   - **Key Point**: The dynamics in multi-agent settings introduce additional complexities, making MARL a rich research area.

7. **Conclusion**
   - As policy gradient methods evolve, the future will focus on efficiency, adaptability, and interpretability in multi-agent contexts.

## Policy Gradient Theorem

The Policy Gradient Theorem provides the foundation for policy optimization:

$$\nabla J(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[ \nabla \log \pi_\theta(a|s) \cdot Q^\pi(s, a) \right] \tag{14}$$

This equation underscores the core of policy optimization and guides future implementations of policy gradient methods.