

# Deep Learning-Policy Gradient Methods

---

## Policy Gradient Methods

---

Policy gradient methods are a type of reinforcement learning algorithm that learn to improve the policy (or mapping from states to actions) by directly optimizing the expected cumulative reward. This is in contrast to value-based methods, which estimate the state-action values and then use these estimates to improve the policy.

### Key Components

1. **Policy:** The function that maps states to actions.
2. **Reward Function:** Defines the feedback received after each action.
3. **Value Function:** Estimates the expected cumulative reward for a given state (optional).
4. **Exploration-Exploitation Tradeoff:** Balances between exploring new actions and exploiting known good policies.

### Policy Gradient Algorithms

#### 1. Reinforce Algorithm

- Updates policy using Monte Carlo estimates of return.
- Pros: Simple, easy to implement.
- Cons: Requires large number of samples for convergence.

#### 2. Proximal Policy Optimization (PPO)

- Uses trust region optimization to constrain updates.
- Pros: More stable and efficient than vanilla policy gradient.
- Cons: Requires careful tuning of hyperparameters.

### Example in PyTorch

Let's consider a simple example of a cart-pole environment, where the goal is to balance the pole upright by controlling the cart's velocity. We'll use the REINFORCE algorithm for illustration purposes.

```

import torch
import gym

# Define the policy network (simple neural net)
class PolicyNetwork(torch.nn.Module):
    def __init__(self, input_dim, output_dim):
        super(PolicyNetwork, self).__init__()
        self.fc1 = torch.nn.Linear(input_dim, 64)
        self.fc2 = torch.nn.Linear(64, output_dim)

    def forward(self, x):
        x = torch.relu(self.fc1(x))
        return torch.tanh(self.fc2(x))

# Initialize the policy network and optimizer
policy_net = PolicyNetwork(input_dim=4, output_dim=1) # cart-pole has 4 state dimensions
optimizer = torch.optim.Adam(policy_net.parameters(), lr=0.01)

# Define the REINFORCE algorithm
def reinforce_algorithm(env, policy_net, optimizer):
    episode_rewards = []
    for _ in range(100): # number of episodes to run
        state = env.reset()
        done = False
        rewards = 0

        while not done:
            action = policy_net(torch.tensor(state)).item() # select an action based on current state
            next_state, reward, done, _ = env.step(action)
            rewards += reward
            state = next_state

        episode_rewards.append(rewards)

    optimizer.zero_grad()

    # Compute policy gradient using Monte Carlo estimates of return
    returns = torch.tensor(episode_rewards).mean()
    loss = -returns * torch.sum(policy_net(torch.tensor(env.observation_space.sample())))

    loss.backward()
    optimizer.step()

# Run the REINFORCE algorithm for multiple episodes
for i in range(1000): # number of iterations to run
    reinforce_algorithm(gym.make('CartPole-v1'), policy_net, optimizer)

```

In this example, we define a simple neural network as our policy function and use the REINFORCE algorithm to update it based on Monte Carlo estimates of return. We iterate for multiple episodes, updating the policy parameters after each episode.

Note that this is a highly simplified example and in practice you would likely want to use more sophisticated algorithms like PPO or trust region optimization to improve stability and efficiency. Additionally, you may need to tune hyperparameters such as learning rates, batch sizes, and exploration-exploitation tradeoffs for optimal performance.

---

*Curated by Brajesh Kumar*