



# AIL 722: Reinforcement Learning

## Lecture 33: Bias and Variance

Raunak Bhattacharyya



**ScAI**

YARDI SCHOOL OF ARTIFICIAL INTELLIGENCE  
INDIAN INSTITUTE OF TECHNOLOGY DELHI

# Recap & Today's Outline

- Policy gradient methods
- Gradient expression
- Reinforce algorithm
- Bias and variance
- Variance reduction
- Baseline

# The Reinforce Algorithm

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \right) \left( \sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right)$$

## REINFORCE: Monte-Carlo Policy-Gradient Control (episodic) for $\pi_*$

Input: a differentiable policy parameterization  $\pi(a|s, \theta)$

Algorithm parameter: step size  $\alpha > 0$

Initialize policy parameter  $\theta \in \mathbb{R}^{d'}$  (e.g., to  $\mathbf{0}$ )

Loop forever (for each episode):

    Generate an episode  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ , following  $\pi(\cdot|\cdot, \theta)$

    Loop for each step of the episode  $t = 0, 1, \dots, T - 1$ :

$$G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k \tag{G_t}$$

$$\theta \leftarrow \theta + \alpha \gamma^t G \nabla \ln \pi(A_t | S_t, \theta)$$

# Gradient Estimator

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \right) \left( \sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right)$$

$$\widehat{\nabla_{\theta} J(\theta)} = \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \right) \left( \sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right)$$

A sample mean of the quantity  $\nabla_{\theta} \log \pi_{\theta}(\tau_i) r(\tau_i)$

**From Sutton and Barto: “...since a Monte-Carlo estimator, will suffer from high variance.”**

$$\mathbb{E}[f(x)] = \frac{1}{n} \sum_{i=1}^N f(x_i)$$

$$\text{Var}(\mathbb{E}[f(x)]) = \frac{\text{Var}(f(x))}{N}$$

**Variance of the estimator is determined by variance of f**

# Why Does Variance Matter?

