



# AIL 722: Reinforcement Learning

## Lecture 35: Variance Reduction

Raunak Bhattacharyya



**ScAI**

YARDI SCHOOL OF ARTIFICIAL INTELLIGENCE  
INDIAN INSTITUTE OF TECHNOLOGY DELHI

# Recap & Today's Outline

- Policy gradients
- Bias and variance
- Baseline
- Optimal baseline
- Causality
- Actor-critic

# Baseline

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

$$\text{Let } Y = \nabla_{\theta} \log p_{\theta}(\tau) \cdot (r(\tau) - b)$$

$$\mathbb{E}[Y] = \mathbb{E}[\nabla_{\theta} \log p_{\theta}(\tau) \cdot r(\tau)] - b \cdot \mathbb{E}[\nabla_{\theta} \log p_{\theta}(\tau)]$$

$$E[\nabla_{\theta} \log p_{\theta}(\tau)] = \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau) d\tau$$

$$= \int \nabla_{\theta} p_{\theta}(\tau) d\tau$$

$$= \nabla_{\theta} \int p_{\theta}(\tau) d\tau = \nabla_{\theta} 1$$

$$E[Y] = E[\nabla_{\theta} \log p_{\theta}(\tau) \cdot r(\tau)]$$

**Our revised estimator Y does not introduce bias**

# Optimal Baseline

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

$$Y = \nabla_{\theta} \log p_{\theta}(\tau) \cdot (r(\tau) - b)$$

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

**Thus, sufficient to analyse variance of Y. Then divide by N**

$$\text{Var}(Y) = \mathbb{E}[Y^2] - (\mathbb{E}[Y])^2$$

**Can we ignore the second term?**

$$E[Y] = E\left[\nabla_{\theta} \log p_{\theta}(\tau) \cdot r(\tau)\right]$$

**Independent of b**

Therefore, we are interested in  $\mathbb{E}[Y^2]$

# Optimal Baseline

$$Y = \nabla_{\theta} \log p_{\theta}(\tau) \cdot (r(\tau) - b)$$

**Goal: Find  $b$  such that the variance is lowest**

Therefore, we are interested in  $\mathbb{E}[Y^2]$

$$\mathbb{E} \left[ g(\tau)^2 (r(\tau) - b)^2 \right]$$

$$\begin{aligned} \frac{d\text{Var}}{db} &= \frac{d}{db} \mathbb{E} [g(\tau)^2 (r(\tau) - b)^2] \\ &= \frac{d}{db} \mathbb{E} [g(\tau)^2 r(\tau)^2] - 2\mathbb{E} [g(\tau)^2 r(\tau) b] + b^2 \mathbb{E} [g(\tau)^2] \\ &= -2\mathbb{E} [g(\tau)^2 r(\tau)] + 2b \mathbb{E} [g(\tau)^2] \\ &= 0 \end{aligned}$$

$$b^{opt} = \frac{\mathbb{E} [g(\tau)^2 r(\tau)]}{\mathbb{E} [g(\tau)^2]}$$

**In practice, we use the average reward as the baseline**

# From Sutton and Barto

## REINFORCE with Baseline (episodic), for estimating $\pi_{\theta} \approx \pi_*$

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

Input: a differentiable state-value function parameterization  $\hat{v}(s, \mathbf{w})$

Algorithm parameters: step sizes  $\alpha^{\theta} > 0$ ,  $\alpha^{\mathbf{w}} > 0$

Initialize policy parameter  $\theta \in \mathbb{R}^{d'}$  and state-value weights  $\mathbf{w} \in \mathbb{R}^d$  (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  $t = 0, 1, \dots, T - 1$ :

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

$$\delta \leftarrow G - \hat{v}(S_t, \mathbf{w})$$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha^{\mathbf{w}} \delta \nabla \hat{v}(S_t, \mathbf{w})$$

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

# Variance Reduction: Causality

$$\begin{aligned}\nabla_{\theta} J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \left( \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \right) \left( \sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right) \\ &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \left( \sum_{t'=1}^T r(s_{i,t'}, a_{i,t'}) \right) \\ &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \left( \sum_{t'=t}^T r(s_{i,t'}, a_{i,t'}) \right) \\ &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) \hat{Q}_{i,t}\end{aligned}$$

Reward-to-go

Reduced variance since we reduced the sum to be a smaller number