



Computer Engineering Department

Reinforcement Learning: Advanced Policy Gradients

Mohammad Hossein Rohban, Ph.D.

Spring 2025

Courtesy: Most of slides are adopted from CS 285 Berkeley.

Lecture 10 - 1

Model-Free RL

- Value-based methods
 - Learnt value function
 - Implicit policy
- Policy-based methods
 - No value function
 - Learnt policy
- Actor-critic methods
 - Learnt value function
 - Learnt policy



Overview of Modern RL Methods



Policy Gradient Intuition

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left(\sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^{T} r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) \right)$$
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \underbrace{\nabla_{\theta} \log \pi_{\theta}(\tau_{i}) r(\tau_{i})}_{T} \qquad \text{maximum likelihood:} \quad \nabla_{\theta} J_{\mathrm{ML}}(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \nabla_{\theta} \log \pi_{\theta}(\tau_{i})$$
$$\sum_{t=1}^{T} \nabla_{\theta} \log_{\theta} \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t})$$

- Good stuff is made more likely
- Bad stuff is made less likely
- Simply formalizes the notion of "trial and error"!



Bias and Variance of Policy Gradient

Unbiased estimation:

$$E\left[\frac{1}{N}\sum_{i=1}^{N}\nabla_{\theta}\log\pi_{\theta}(\tau_{i})r(\tau_{i})\right] = \nabla_{\theta}J(\theta)$$

But suffers from high variance!

Reducing Variance

- Everything in the gradient whose expected is zero could be removed, without affecting the optimization, but could lead to lower gradient variance!
- Causality trick
- Discount factor
- Baseline
- Actor-critic
- Optimization techniques:
 - Natural gradient
 - Trust region

Preducing Variance: Discount Factor

$$r(s, a, t) + r(s, a, t) + r(s, a, t) + r(s, t) = r(s, t)$$
option 1: $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\sum_{t'=t}^{T} \gamma^{t'-t} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right)$
Not the same
option 2: $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \left(\sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \right) \left(\sum_{t=1}^{T} \gamma^{t-1} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right)$

$$Iff \left(\nabla J - Z \right) = Iff \left(\nabla J \right)$$
 $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\sum_{t=1}^{T} \gamma^{t-1} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right)$
 $Ff(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\sum_{t'=t}^{T} \gamma^{t'-t} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right)$
 $Ff(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\sum_{t'=t}^{T} \gamma^{t'-t} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right)$

Vor [ال الحري (ح) [r(ح) - ل] Reducing Variance: Baselines

$$T = \left(\begin{array}{c} S_{\ldots} a_{\ldots} S_{1} a_{1} \cdots \right) \qquad \text{a convenient identity} \\ p_{\theta}(\tau) \otimes \frac{1}{N} \sum_{i=1}^{N} \nabla_{\theta} \log p_{\theta}(\tau) [r(\tau) - b] \qquad \text{a convenient identity} \\ b = \frac{1}{N} \sum_{i=1}^{N} r(\tau) \qquad \text{idenly:} \qquad b = IE \left[\begin{array}{c} r(\tau) \\ \tau \sim \pi_{\theta} \end{array} \right] \\ V_{\theta r} \left[\begin{array}{c} X_{1} - X_{2} \end{array} \right] \\ V_{\theta r} \left[\begin{array}{c} X_{1} - X_{2} \end{array} \right] \\ E[\nabla_{\theta} \log p_{\theta}(\tau)b] = \int p_{\theta}(\tau) \nabla_{\theta} \log p_{\theta}(\tau)b \, d\tau = \int \nabla_{\theta} p_{\theta}(\tau)b \, d\tau = b \nabla_{\theta} \int p_{\theta}(\tau) d\tau = b \nabla_{\theta} 1 = 0 \\ \text{subtracting a baseline is unbiased in expectation!} \\ \text{average reward is not the best baseline, but it's pretty good!} \qquad IE \left[\left(\begin{array}{c} X_{1} - X_{2} - \begin{array}{c} V_{x_{1}} + \begin{array}{c} V_{x_{2}} \end{array} \right)^{2} \right] \\ = \sqrt{m} \left[X_{1} \right] + V_{xr} \left[X_{2} \right] - 2C \text{ ov } \cdots \end{array} \right]$$

subtracting a baseline is *unbiased* in expectation!

average reward is *not* the best baseline, but it's pretty good!

Reducing Variance: Baselines

Faster convergence:



Reducing Variance: Review

- Exploiting causality
 - Future doesn't affect the past
- Discount factor
 - Two different version
- Baselines
 - Analyzing variance for deriving optimal baselines
- Now: Introducing actor-critic methods!

Policy Gradients so Far

$\hat{Q}^{\pi}(\mathbf{x}_t, \mathbf{u}_t) = \sum r(\mathbf{x}_{t'}, \mathbf{u}_{t'})$ **REINFORCE** algorithm: 1. sample $\{\tau^i\}$ from $\pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)$ (run the policy) 2. $\nabla_{\theta} J(\theta) \approx \sum_{i} \left(\sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{t}^{i} | \mathbf{s}_{t}^{i}) \left(\sum_{t'=t}^{T} r(\mathbf{s}_{t'}^{i}, \mathbf{a}_{t'}^{i}) \right) \right)$ 3. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$ generate samples (i.e. run the policy $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \hat{Q}_{i,t}^{\pi}$ $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$ "reward to go"

t' = t

fit a model to

estimate return

improve the policy

Improving Estimation of Reward to Go

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(\sum_{t'=1}^{T} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right)$$

 $\hat{Q}_{i,t}$: estimate of expected reward if we take action $\mathbf{a}_{i,t}$ in state $\mathbf{s}_{i,t}$

How to make a better estimate?

$$Q(\mathbf{s}_{t}, \mathbf{a}_{t}) = \sum_{t'=t}^{T} E_{\pi_{\theta}} \left[r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_{t}, \mathbf{a}_{t} \right]: \text{ true } expected \text{ reward-to-go}$$
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) Q(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$



Improving Estimation of Reward to Go

Further improvement: Adding a baseline!

$$Q(\mathbf{s}_{t}, \mathbf{a}_{t}) = \sum_{t'=t}^{T} E_{\pi_{\theta}} \left[r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_{t}, \mathbf{a}_{t} \right]: \text{ true } expected \text{ reward-to-go}$$
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) \left(Q(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) - b_{t} \right)$$

$$b_t = \frac{1}{N} \sum_i Q(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$





Lecture 10 - 15

Advantage Value

 $Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t) = \sum_{t'=t}^{T} E_{\pi_{\theta}} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t, \mathbf{a}_t]$: total reward from taking \mathbf{a}_t in \mathbf{s}_t

 $V^{\pi}(\mathbf{s}_t) = E_{\mathbf{a}_t \sim \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t)}[Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t)]$: total reward from \mathbf{s}_t

 $A^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) = Q^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) - V^{\pi}(\mathbf{s}_{t}): \text{ how much better } \mathbf{a}_{t} \text{ is}$ $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) A^{\pi}(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$

the better this estimate, the lower the variance

Advantage Value Approximation

$$Q^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) = \sum_{t'=t}^{T} E_{\pi_{\theta}} \left[r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_{t}, \mathbf{a}_{t} \right]$$
$$V^{\pi}(\mathbf{s}_{t}) = E_{\mathbf{a}_{t} \sim \pi_{\theta}(\mathbf{a}_{t} | \mathbf{s}_{t})} \left[Q^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) \right]$$
$$A^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) = Q^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) - V^{\pi}(\mathbf{s}_{t})$$
$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_{i,t} | \mathbf{s}_{i,t}) A^{\pi}(\mathbf{s}_{i,t}, \mathbf{a}_{i,t})$$

fit
$$Q^{\pi}$$
, V^{π} , or A^{π}
fit a model to
estimate return
generate
samples (i.e.
run the policy)
 $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

$$Q^{\pi}(\mathbf{s}_{t}, \mathbf{a}_{t}) = \sum_{t'=t}^{T} E_{\pi_{\theta}} \left[r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_{t}, \mathbf{a}_{t} \right]$$
$$= r(\mathbf{s}_{t}, \mathbf{a}_{t}) + \sum_{t'=t+1}^{T} E_{\pi_{\theta}} \left[r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_{t}, \mathbf{a}_{t} \right]$$

 $\approx V^{\pi}(\mathbf{s}_{t+1})$

Advantage Value Approximation

 $Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t) \approx r(\mathbf{s}_t, \mathbf{a}_t) + V^{\pi}(\mathbf{s}_{t+1})$

 $A^{\pi}(\mathbf{s}_t, \mathbf{a}_t) \approx r(\mathbf{s}_t, \mathbf{a}_t) + V^{\pi}(\mathbf{s}_{t+1}) - V^{\pi}(\mathbf{s}_t)$



 \mathbf{S}

let's just fit $V^{\pi}(\mathbf{s})!$

Lecture 10 - 18

$$V^{\pi}(\mathbf{s}_t) = \sum_{t'=t}^T E_{\pi_{\theta}} \left[r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t \right]$$
$$J(\theta) = E_{\mathbf{s}_1 \sim p(\mathbf{s}_1)} \left[V^{\pi}(\mathbf{s}_1) \right]$$

how can we perform policy evaluation?

Monte Carlo policy evaluation (this is what policy gradient does)

$$V^{\pi}(\mathbf{s}_t) \approx \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'})$$
$$V^{\pi}(\mathbf{s}_t) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'})$$

(requires us to reset the simulator)



Monte Carlo estimation with function approximator:

 $V^{\pi}(\mathbf{s}_t) \approx \sum_{t'=t}^T r(\mathbf{s}_{t'}, \mathbf{a}_{t'})$

not as good as this: $V^{\pi}(\mathbf{s}_t) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t'=t}^{T} r(\mathbf{s}_{t'}, \mathbf{a}_{t'})$

but still pretty good!

training data: $\left\{ \left(\mathbf{s}_{i,t}, \sum_{t'=t}^{T} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'}) \right) \right\}$

supervised regression:
$$\mathcal{L}(\phi) = \frac{1}{2} \sum_{i} \left\| \hat{V}_{\phi}^{\pi}(\mathbf{s}_{i}) - y_{i} \right\|^{2}$$





How to make a better estimate?

ideal target:
$$y_{i,t} = \sum_{t'=t}^{T} E_{\pi_{\theta}} \left[r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_{i,t} \right] \approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + V^{\pi}(\mathbf{s}_{i,t+1}) \approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \hat{V}^{\pi}_{\phi}(\mathbf{s}_{i,t+1})$$

Monte Carlo target: $y_{i,t} = \sum_{t'=t}^{T} r(\mathbf{s}_{i,t'}, \mathbf{a}_{i,t'})$

directly use previous fitted value function!

Bootstrap Estimation with Function Approximator

ideal target:
$$y_{i,t} = \sum_{t'=t}^{T} E_{\pi_{\theta}} [r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_{i,t}] \approx r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \hat{V}_{\phi}^{\pi}(\mathbf{s}_{i,t+1})$$

training data:
$$\left\{ \left(\mathbf{s}_{i,t}, r(\mathbf{s}_{i,t}, \mathbf{a}_{i,t}) + \hat{V}^{\pi}_{\phi}(\mathbf{s}_{i,t+1}) \right) \right\}$$

supervised regression:
$$\mathcal{L}(\phi) = \frac{1}{2} \sum_{i} \left\| \hat{V}_{\phi}^{\pi}(\mathbf{s}_{i}) - y_{i} \right\|^{2}$$

Batch Actor-Critic Algorithm

batch actor-critic algorithm:





 $V^{\pi}(\mathbf{s}_t) = \sum_{t'=t}^{T} E_{\pi_{\theta}} \left[r(\mathbf{s}_{t'}, \mathbf{a}_{t'}) | \mathbf{s}_t \right]$

Actor-Critic Algorithm: Architecture Design

batch actor-critic algorithm:



+ simple & stable

- no shared features between actor & critic



Actor-Critic Algorithm: Batch vs. Online

batch actor-critic algorithm:

1. sample
$$\{\mathbf{s}_i, \mathbf{a}_i\}$$
 from $\pi_{\theta}(\mathbf{a}|\mathbf{s})$ (run it on the robot)
2. fit $\hat{V}^{\pi}_{\phi}(\mathbf{s})$ to sampled reward sums
3. evaluate $\hat{A}^{\pi}(\mathbf{s}_i, \mathbf{a}_i) = r(\mathbf{s}_i, \mathbf{a}_i) + \gamma \hat{V}^{\pi}_{\phi}(\mathbf{s}'_i) - \hat{V}^{\pi}_{\phi}(\mathbf{s}_i)$
4. $\nabla_{\theta} J(\theta) \approx \sum_i \nabla_{\theta} \log \pi_{\theta}(\mathbf{a}_i|\mathbf{s}_i) \hat{A}^{\pi}(\mathbf{s}_i, \mathbf{a}_i)$
5. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

online actor-critic algorithm:

Proximal Policy Optimization

 We don't want the new policy to change a lot in an iteration. Why?

$D_{KL}(\pi_{\theta'}(.|s) || \pi_{\theta}(.|s)) \le \epsilon$

- What is the effect of the constraint?
- Recall KL-Divergence:

$$D_{KL}(\pi_{\theta'}(.|s) | | \pi_{\theta}(.|s)) = \Sigma_a \pi_{\theta'}(a|s) \log \frac{\pi_{\theta'}(a|s)}{\pi_{\theta}(a|s)}$$

We are effectively constraining the ratio $\frac{\pi_{\theta'}(a|s)}{\pi_{\theta}(a|s)}$