

ECE/SIOC 228 Machine Learning for Physical Applications

Lecture 15: Deep Reinforcement Learning II

Yuanyuan Shi

Assistant Professor, ECE

University of California, San Diego

Review: Policy Gradient

UC San Diego

$$\text{Episodic task: } J(\theta) = E_{z \sim p_\theta(z)} \left[\underbrace{\sum_{t=0}^{T-1} \gamma^t R(s_t, a_t)}_{\gamma(z)} \right]$$

$\gamma(z)$ is the return (total discounted reward) of an episode

$$z \doteq s_0, a_0, s_1, a_1, s_2, a_2, \dots, s_{T-1}, a_{T-1}, s_T$$

Find the best θ that optimize $J(\theta)$:

$$\theta^* = \underset{\theta}{\operatorname{argmax}} J(\theta) = \underset{\theta}{\operatorname{argmax}} E_{z \sim p_\theta(z)} [\gamma(z)]$$

Review: Policy Gradient

UC San Diego

- Use Gradient Ascent. $\theta_{t+1} \leftarrow \theta_t + \alpha \nabla_{\theta} J(\theta)$
- $\nabla_{\theta} J(\theta) = \nabla_{\theta} E_z[r(z)] = \nabla_{\theta} \int p_{\theta}(z) r(z) dz = \int \nabla_{\theta} p_{\theta}(z) r(z) dz$
 $= \int p_{\theta}(z) \frac{\nabla_{\theta} p_{\theta}(z)}{p_{\theta}(z)} r(z) dz = E_{z \sim p_{\theta}(z)} [\nabla_{\theta} \log p_{\theta}(z) r(z)]$
 $= E_{z \sim p_{\theta}(z)} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) r(z) \right]$

A property used in derivation; $p_{\theta}(z) = p_{\theta}(s_0) \prod_{t=0}^{T-1} \pi_{\theta}(a_t | s_t) p(s_{t+1} | s_t, a_t)$

Note:

- In some reference, time index starts from $t=1, \dots, T$
- If T is finite, we can also choose $\gamma \in (0, 1]$

The REINFORCE Algorithm

UC San Diego

Standard Policy gradient:

$$\nabla_{\theta} J(\theta) = E_{z \sim p_{\theta}(z)} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) r(z) \right]$$

$$\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \underbrace{\sum_{t=t}^{T-1} y^t r(s_t^i, a_t^i)}$$

* Causality: policy at time t^* cannot affect reward at time $t < t^*$

* Fix: $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) \sum_{t=t}^{T-1} \underbrace{y^{t-t}}_{\text{return: } G_t^i} r(s_t^i, a_t^i)$

REINFORCE Algorithm

1. sample $\{\tau^i\}$ from $\pi_{\theta}(a_t | s_t)$ (run the policy)
2. $\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_i \sum_t \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) G_t^i$
3. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

score function.

return: G_t^i //

Example: Gaussian Policies

UC San Diego

- For continuous action space, a Gaussian policy is natural

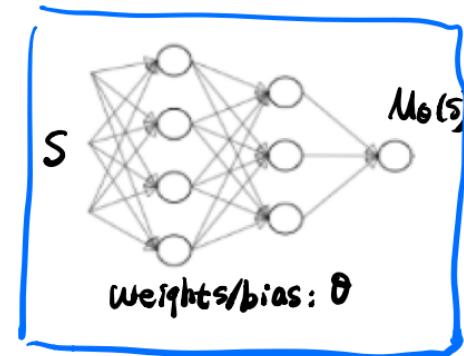
- Gaussian policy: $a \sim N(\mu_\theta(s), \sigma^2)$

$$\pi_\theta(a|s) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{a-\mu_\theta(s)}{\sigma}\right)^2\right]$$

$$\log \pi_\theta(a|s) = -\log \sigma\sqrt{2\pi} - \frac{1}{2}\left(\frac{a-\mu_\theta(s)}{\sigma}\right)^2$$

- Score function:

$$\nabla_\theta \log \pi_\theta(a|s) = \frac{(a-\mu_\theta(s)) \nabla_\theta \mu_\theta(s)}{\sigma^2}$$



Example: Softmax Policies

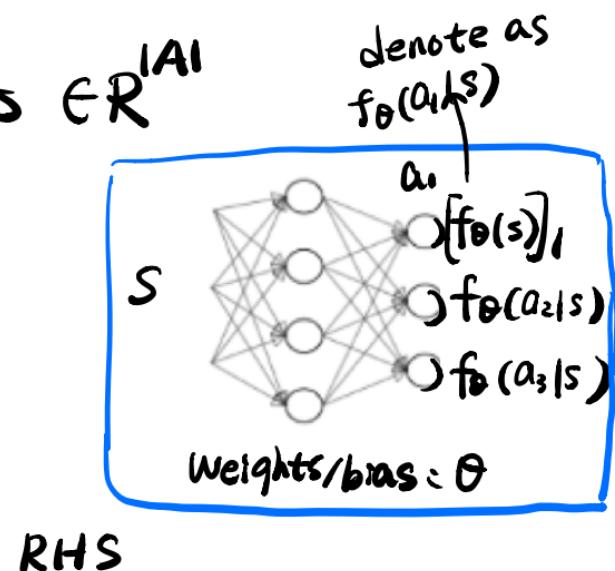
UC San Diego

- For discrete action space, we commonly use a **categorical policy** (use softmax activation)

$$\text{Softmax}(z) = \frac{e^z}{\sum_{k=1}^{|A|} e^{z_k}} \leftarrow \text{probabilities } \in \mathbb{R}^{|A|}$$

$$\bullet \text{ For example, } \pi_\theta(a_i|s) = \frac{e^{f_\theta(a_i|s)}}{\sum_{k=1}^{|A|} e^{f_\theta(a_k|s)}}$$

$$\log \pi_\theta(a_i|s) = f_\theta(a_i|s) - \underbrace{\log \sum_{k=1}^{|A|} e^{f_\theta(a_k|s)}}_{\text{RHS}}$$



Example: Softmax Policies

UC San Diego

Score function: $\nabla_{\theta} \log \pi_{\theta}(a_i|s) = \nabla_{\theta} f_{\theta}(a_i|s) - \nabla_{\theta} \text{RHS}$

$$\nabla_{\theta} \text{RHS} = \nabla_{\theta} \log \sum_{k=1}^{|A|} e^{f_{\theta}(a_k|s)}$$

$$= \frac{\nabla_{\theta} \sum_{k=1}^{|A|} e^{f_{\theta}(a_k|s)}}{\sum_{k=1}^{|A|} e^{f_{\theta}(a_k|s)}} = \frac{\sum_{k=1}^{|A|} e^{f_{\theta}(a_k|s)} \cdot \nabla_{\theta} f_{\theta}(a_k|s)}{\sum_{k=1}^{|A|} e^{f_{\theta}(a_k|s)}}$$

$$= \sum_{k=1}^{|A|} \pi_{\theta}(a_k|s) \nabla_{\theta} f_{\theta}(a_k|s)$$

$$\nabla \log \pi_{\theta}(a_i|s) = \nabla_{\theta} f_{\theta}(a_i|s) - E_{a \sim \pi_{\theta}(\cdot|s)} [\nabla_{\theta} f_{\theta}(a|s)]$$

Implementing Policy Gradients

UC San Diego

Pseudocode example:

REINFORCE Algorithm:

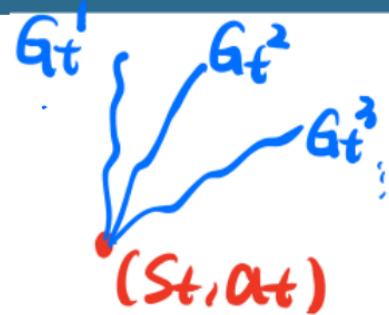
```
# Given:  
# states: [N*T, state_dim] tensor of states  
# actions: [N*T, action_dim] tensor of actions  
# returns: [N*T, 1] tensor of returns  
  
update_policy(states, actions, returns):  
    logprob = calculate_logprob(states, actions) #this should return [N*T, 1] actions log probability  
    loss = -torch.mean(logprob*returns)  
    optimizer.zero_grad()  
    loss.backward()  
    optimizer.step()
```

$\log \pi_\theta(a_t^i | s_t^i)$

$\uparrow \log \pi_\theta(a_t^i | s_t^i) \cdot G_t^i$

REINFORCE Algorithm

$$\begin{aligned}\nabla_{\theta} J(\theta) &= E_{z \sim p_{\theta}(z)} \left[\sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) G_t \right] \\ &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t^i | s_t^i) G_t^i\end{aligned}$$

REINFORCE with Baseline

$$\nabla J(\theta) = E_{z \sim p_{\theta}(z)} \left[\sum_{t=0}^{T-1} \nabla \log \pi_{\theta}(a_t | s_t) (G_t - b(s_t)) \right]$$

To see why, we will show

$$E_{z \sim p_{\theta}(z)} \left[\sum_{t=0}^{T-1} \nabla \log \pi_{\theta}(a_t | s_t) b(s_t) \right] = 0$$

Can we do that??

REINFORCE with Baseline

UC San Diego

Due to linearity of expectation, it suffices to show that for any time step t , $E_{\zeta \sim P_\theta(z)} [\nabla_\theta \log \pi_\theta(a_t | s_t) b(s_t)] = 0$

$$E_{\zeta \sim P_\theta(z)} [\nabla_\theta \log \pi_\theta(a_t | s_t) b(s_t)] \quad s_0, a_0, s_1, a_1, \dots, s_t \mid a_t, \dots, s_{t-1}, a_{t-1}, s_t$$

$$= E_{s_0:t, a_0:t-1} [E_{s_{t+1:T}, a_{t:T}} [\nabla_\theta \log \pi_\theta(a_t | s_t) b(s_t)]]$$

$$= E_{s_0:t, a_0:t-1} [b(s_t) E_{s_{t+1:T}, a_{t:T}} [\nabla_\theta \log \pi_\theta(a_t | s_t)]]$$

$$= E_{s_0:t, a_0:t-1} [b(s_t) \underbrace{E_{a_t} [\nabla_\theta \log \pi_\theta(a_t | s_t)]}_{\text{blue underline}}]$$

REINFORCE with Baseline

UC San Diego

$$\begin{aligned}
 E_{\pi_\theta} [\nabla_\theta \log \pi_\theta (a_t | s_t)] &= \int \pi_\theta(a_t | s_t) \frac{\nabla_\theta \pi_\theta (a_t | s_t)}{\pi_\theta (a_t | s_t)} da_t \\
 &= \nabla_\theta \int \pi_\theta (a_t | s_t) da_t \\
 &= \nabla_\theta 1 = 0
 \end{aligned}$$

Thus: $E_{z \sim p_\theta(z)} [\nabla_\theta \log \pi_\theta (a_t | s_t) b(s_t)] = 0, \forall t$

~~Subtracting a baseline $b(s_t)$ is unbiased in expectation~~

$$E_{z \sim p_\theta(z)} \left[\sum_{t=0}^{T-1} \nabla \log \pi_\theta (a_t | s_t) (G_t - b(s_t)) \right] = E_{z \sim p_\theta(z)} \left[\sum_{t=0}^{T-1} \nabla \log \pi_\theta (a_t | s_t) G_t \right]$$

Then ... Why subtract baseline $b(s_t)$? \Rightarrow reduce variance

$$\text{Var}\left(\sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta(a_t | s_t) (G_t - b(s_t))\right)$$

$$\approx \sum_{t=0}^{T-1} \text{Var}\left(\nabla_\theta \log \pi_\theta(a_t | s_t) (G_t - b(s_t))\right) \quad (\text{assume independence})$$

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

$$= \sum_{t=0}^{T-1} \underbrace{\mathbb{E}\left[\left(\nabla_\theta \log \pi_\theta(a_t | s_t) (G_t - b(s_t))\right)^2\right]}_{①} - \underbrace{\mathbb{E}\left[\nabla_\theta \log \pi_\theta(a_t | s_t) (G_t - b(s_t))\right]}_{②}^2$$

REINFORCE with Baseline

UC San Diego

$$\textcircled{2} = E \left[\nabla_\theta \log \pi_\theta(s_t | a_t) G_t \right]^2$$

$$\textcircled{1} = E \left[(\nabla_\theta \log \pi_\theta(a_t | s_t))^2 (G_t - b(s_t))^2 \right]$$

$$= E \left[(\nabla_\theta \log \pi_\theta(a_t | s_t))^2 (G_t^2 - 2G_t b(s_t) + b(s_t)^2) \right]$$

$$= \underbrace{E \left[(\nabla_\theta \log \pi_\theta(a_t | s_t))^2 G_t^2 \right]}_{\textcircled{1a}} - 2 \underbrace{E \left[(\nabla_\theta \log \pi_\theta(a_t | s_t))^2 G_t b(s_t) \right]}_{\textcircled{1b}}$$

$$+ E \left[(\nabla_\theta \log \pi_\theta(a_t | s_t))^2 b^2(s_t) \right] \quad \textcircled{1c}$$

REINFORCE with Baseline

UC San Diego

$$\text{Var}(\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) g_t) = \textcircled{1a} - \textcircled{2} \quad \text{REINFORCE}$$

$$\begin{aligned} \text{Var}(\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) (g_t - b(s_t))) &= \text{REINFORCE with} \\ &= \textcircled{1a} - \underline{\textcircled{1b}} + \textcircled{1c} - \textcircled{2} \\ &\quad \text{Baseline} \end{aligned}$$

Choose $b^*(s_t)$ that minimize $(-\textcircled{1b} + \textcircled{1c})$, thus to obtain variance reduction

REINFORCE with Baseline

UC San Diego

$$\underbrace{-lb + lc}_{\ell} = E_{s_t, t, a_t, t} \left\{ -2b(s_t) E_{a_{t+1}, s_{t+1}, t} \left[(\nabla_\theta \log \pi_\theta(a_t | s_t))^2 G_t \right] + b^2(s_t) E_{a_{t+1}, s_{t+1}, t} \left[(\nabla_\theta \log \pi_\theta(a_t | s_t))^2 \right] \right\}$$

$$\text{Let } \alpha = E_{a_{t+1}, s_{t+1}, t} \left[(\nabla_\theta \log \pi_\theta(a_t | s_t))^2 | s_t \right]$$

$$\beta = E_{a_{t+1}, s_{t+1}, t} \left[(\nabla_\theta \log \pi_\theta(a_t | s_t))^2 G_t | s_t \right]$$

$$\min_{b(s_t)} -2b(s_t) \cdot \beta + \alpha b^2(s_t) \Rightarrow b^*(s_t) = \frac{\beta}{\alpha}$$

REINFORCE with Baseline

UC San Diego

$$b^*(s_t) = \frac{E_{a_{t:T-1}, s_{t+1:T}} [(\nabla_\theta \log \pi_\theta(a_t | s_t))^2 G_t | s_t]}{E_{a_{t:T-1}, s_{t+1:T}} [(\nabla_\theta \log \pi_\theta(a_t | s_t))^2 | s_t]} \quad \begin{array}{l} \text{← expected return from } s_t, \\ \text{Weighted by } [\nabla_\theta \log \pi_\theta(a_t | s_t)]^2 \end{array}$$

- Therefore, if we drop the "weight terms", we obtain a commonly used baseline; $b(s_t) = V^{\pi_\theta}(s_t) = E_{\pi_\theta}[G_t | s_t] = E_{a_{t:T-1}, s_{t+1:T}}[G_t | s_t]$
- With $b^*(s_t)$, $\ell^* = -\frac{\beta^2}{2} < 0 \Rightarrow$ variance reduction!

Value Function Fitting

UC San Diego

$$V^{\pi_\theta}(s_t) = E_{\pi_\theta}[G_t | s_t]$$

$$= E_{a_{t:T-1}, s_{t+1:T}}[G_t | s_t]$$

(sample)

$$\approx \frac{1}{N} \sum_{i=1}^N \sum_{t'=t}^T \gamma^{t'-t} r(s_{t'}^i, a_{t'}^i)$$

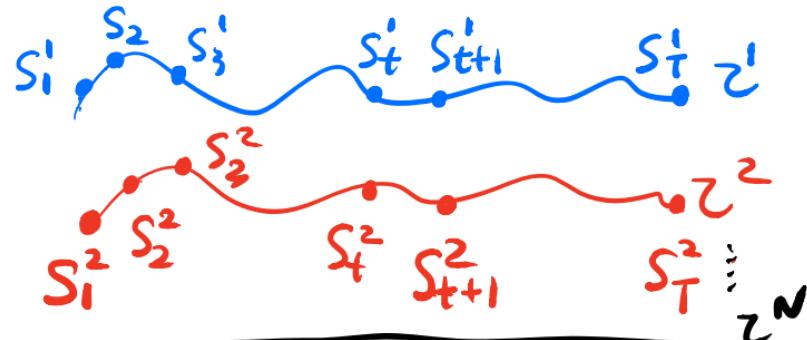


different
episodes
from
 s_t

- Require us to reset simulator to state s_t and generate N episodes.
- If $N=1$, $V^{\pi_\theta}(s_t) \approx \sum_{t'=t}^T \gamma^{t'-t} r(s_{t'}, a_{t'})$ (Can be noisy)

Value Function Fitting

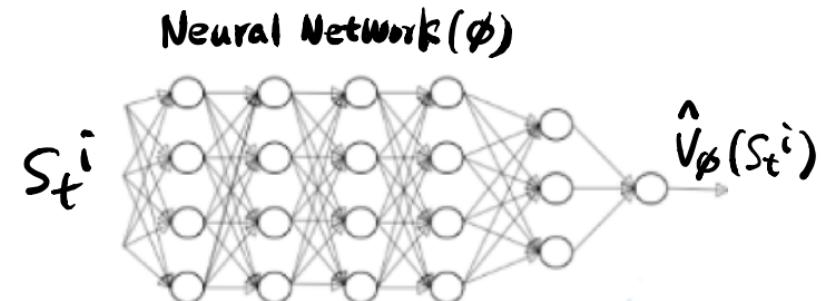
UC San Diego



N episodes with length T

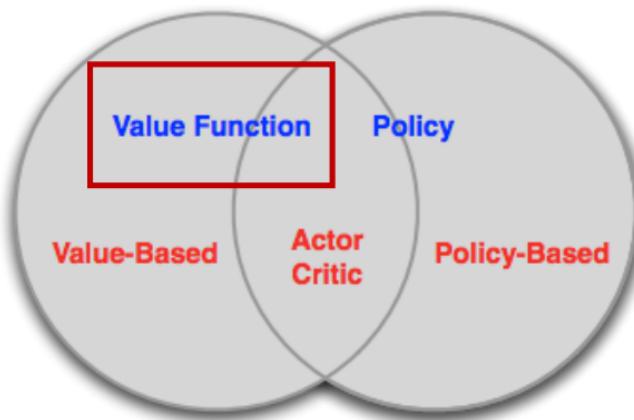
Training data: $\{(S_t^i, \sum_{t'=t}^T \gamma^{t'-t} r(s_{t'}^i, a_{t'}^i))\}$

Supervised learning $L(\phi) = \frac{1}{2} \sum_{i=1}^N \sum_{t=0}^{T-1} (\hat{V}_\phi^{\pi_0}(S_t^i) - G_t^i)^2$



Value Function Fitting

UC San Diego



- More on value-based deep reinforcement learning after 2 guest lectures:
 - Machine Learning for Power Systems
 - Machine Learning for Robotics
- We will also talk about combining policy-based + value-based DRL → **actor-critic method!**