

10703 Deep Reinforcement Learning

Policy Gradient Methods - Part 2

Tom Mitchell

October 3, 2018

Reading: Barto & Sutton, Chapter 13

Used Materials

- Much of the material and slides for this lecture were taken from Chapter 13 of Barto & Sutton textbook.
- Some slides are borrowed from Ruslan Salakhutdinov, who in turn borrowed from Rich Sutton's RL class and David Silver's Deep RL tutorial

Value-Based and Policy-Based RL

▶ Value Based

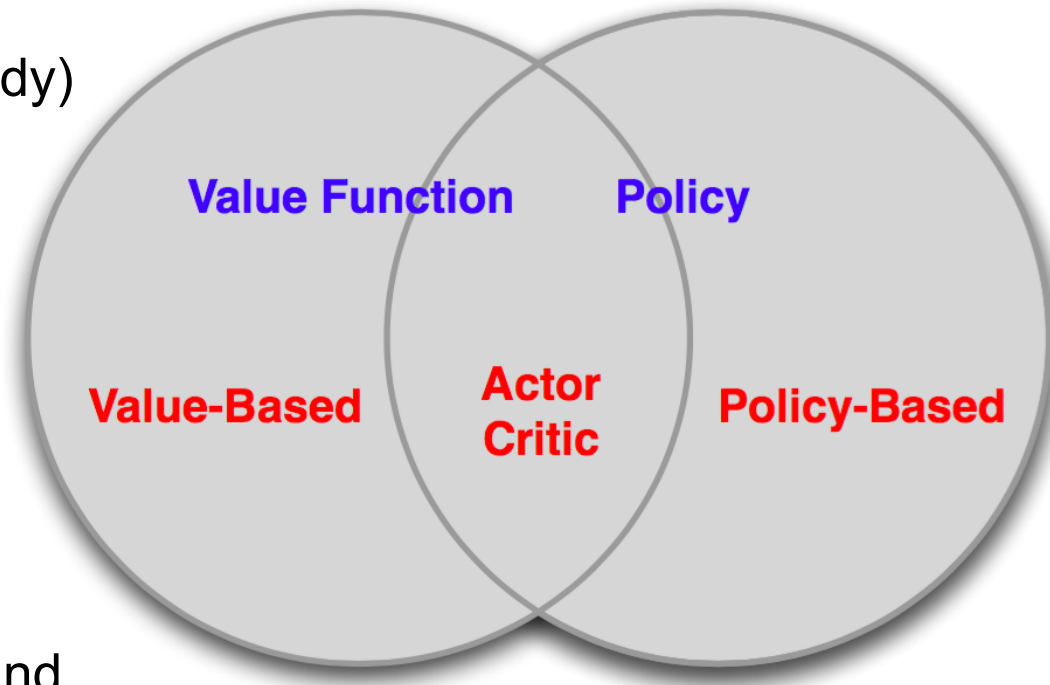
- Learn a Value Function
- Implicit policy (e.g. ϵ -greedy)

▶ Policy Based

- Learn a Policy directly

▶ Actor-Critic

- Learn a Value Function, and
- Learn a Policy



REINFORCE algorithm

REINFORCE: Monte-Carlo Policy-Gradient Control (episodic) for π_*

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

Algorithm parameter: step size $\alpha > 0$

Initialize policy parameter $\boldsymbol{\theta} \in \mathbb{R}^{d'}$ (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 $\pi(\cdot|\cdot, \boldsymbol{\theta})$

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

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

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

Typical Parameterized Differentiable Policy

- ▶ Softmax:

$$\pi(a|s, \boldsymbol{\theta}) \doteq \frac{e^{h(s,a,\boldsymbol{\theta})}}{\sum_b e^{h(s,b,\boldsymbol{\theta})}},$$

where $h(s,a,\boldsymbol{\theta})$ is any function of s , a with params $\boldsymbol{\theta}$
e.g., linear function of features $\mathbf{x}(s,a)$ you make up

$$h(s, a, \boldsymbol{\theta}) = \boldsymbol{\theta}^\top \mathbf{x}(s, a)$$
$$\mathbf{x}(s, a) \in \mathbb{R}^d$$

e.g., $h(s,a,\boldsymbol{\theta})$ is output of trained neural net

Good news:

- REINFORCE converges to local optimum under usual SGD assumptions
- because $E_{\pi}[G_t] = Q(S_t, A_t)$

But variance is high

- recall high variance of Monte Carlo sampling

Two remedies:

- add a baseline (learn diff from baseline)
- Actor-Critic model (learn both Q and π)

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 $t = 0, 1, \dots, T - 1$:

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

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

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

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

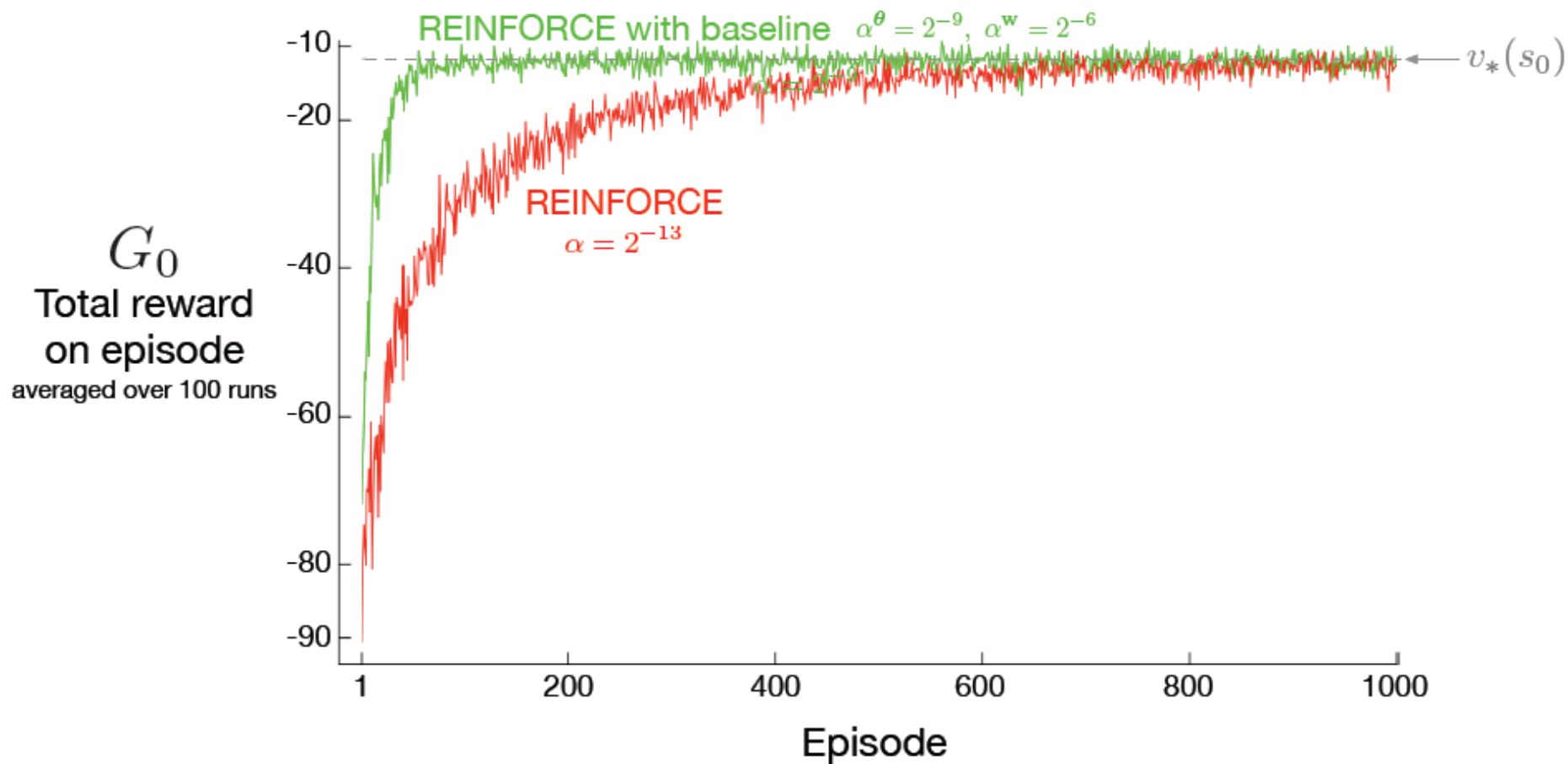
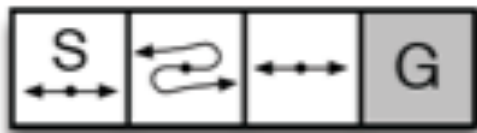


Figure 13.2: Adding a baseline to REINFORCE can make it learn much faster, a

Good news:

- REINFORCE converges to local optimum under usual SGD assumptions
- because $E_{\pi}[G_t] = Q(S_t, A_t)$

But variance is high

- recall high variance of Monte Carlo sampling

Two remedies:

- add a baseline (learn diff from baseline)
- Actor-Critic model (learn both Q and π)

Actor-Critic Model

- learn both Q and π
- use Q to generate target values, instead of G

One step actor-critic model:

$$\begin{aligned}\boldsymbol{\theta}_{t+1} &\doteq \boldsymbol{\theta}_t + \alpha \left(\underline{G_{t:t+1}} - \hat{v}(S_t, \mathbf{w}) \right) \frac{\nabla \pi(A_t | S_t, \boldsymbol{\theta}_t)}{\pi(A_t | S_t, \boldsymbol{\theta}_t)} \\ &= \boldsymbol{\theta}_t + \alpha \left(\underline{R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w})} - \hat{v}(S_t, \mathbf{w}) \right) \frac{\nabla \pi(A_t | S_t, \boldsymbol{\theta}_t)}{\pi(A_t | S_t, \boldsymbol{\theta}_t)} \\ &= \boldsymbol{\theta}_t + \alpha \delta_t \frac{\nabla \pi(A_t | S_t, \boldsymbol{\theta}_t)}{\pi(A_t | S_t, \boldsymbol{\theta}_t)}.\end{aligned}$$

One-step Actor–Critic (episodic), for estimating $\pi_{\theta} \approx \pi_*$

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

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

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):

 Initialize S (first state of episode)

$I \leftarrow 1$

 Loop while S is not terminal (for each time step):

$A \sim \pi(\cdot|S, \theta)$

 Take action A , observe S', R

$\delta \leftarrow R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})$ (if S' is terminal, then $\hat{v}(S', \mathbf{w}) \doteq 0$)

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

$\theta \leftarrow \theta + \alpha^{\theta} I \delta \nabla \ln \pi(A|S, \theta)$

$I \leftarrow \gamma I$

$S \leftarrow S'$

Summary:

Monte Carlo Policy Gradient:

$$\theta_{t+1} \doteq \theta_t + \alpha G_t \frac{\nabla \pi(A_t | S_t, \theta_t)}{\pi(A_t | S_t, \theta_t)}$$

Monte-Carlo Policy Gradient with baseline:

$$\theta_{t+1} \doteq \theta_t + \alpha \left(G_t - b(S_t) \right) \frac{\nabla \pi(A_t | S_t, \theta_t)}{\pi(A_t | S_t, \theta_t)}$$

Actor-Critic Policy Gradient

$$\theta_{t+1} = \theta_t + \alpha \left(R_t + \gamma \hat{v}(S_{t+1}) - \hat{v}(S_t) \right) \frac{\nabla \pi(A_t | S_t, \theta_t)}{\pi(A_t | S_t, \theta_t)}$$