

Approximate Policy Iteration and Policy-Based Learning Methods

Chen-Yu Wei

Approximate Policy Iteration (API)

For $k = 1, 2, \dots$

Evaluate $\hat{Q}_k \approx Q^{\pi_k}$

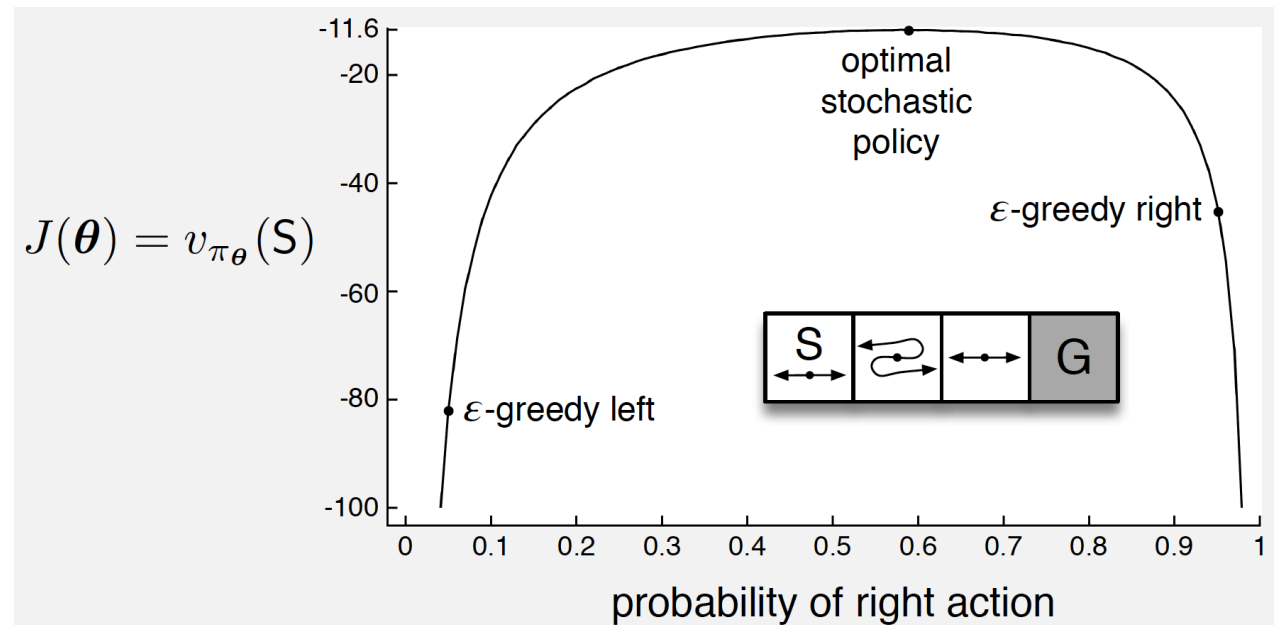
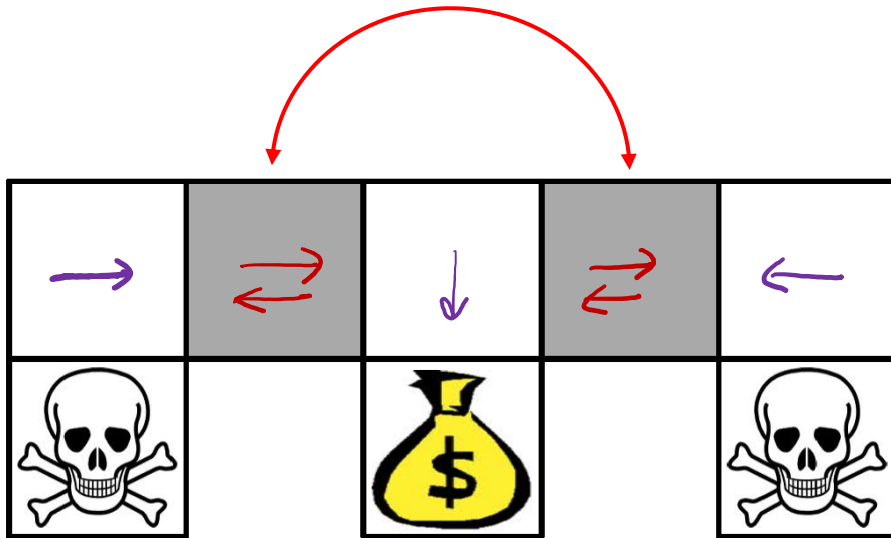
$\pi_{k+1}(s) \leftarrow \operatorname{argmax}_a \hat{Q}_k(s, a)$

Value-based : $\overset{Q^z}{\hat{Q}^*, V^z, V^*} \approx \boxed{V_\theta}$
Policy-based : $\pi_\theta(a/s)$

Limitation (shared by all value-based methods):

- Vulnerable to representation error

Limitation of Value Function Approximation



Incremental Policy Updates through Exponential Weights

For $k = 1, 2, \dots$

Evaluate $\hat{Q}_k \approx Q^{\pi_k}$

Perform incremental policy update such as

$$\pi_{k+1}(a|s) \propto \pi_k(a|s) \exp(\eta \hat{Q}_k(s, a)) \Leftrightarrow \pi_{\theta_{k+1}}(\cdot|s) = \underset{\pi}{\operatorname{argmax}} \left\{ \sum_a \pi(a|s) \hat{Q}_k(s, a) - \frac{1}{2} \operatorname{KL}(\pi(\cdot|s), \pi_{\theta_k}(\cdot|s)) \right\}$$

Function approximation for policies:

$$\theta_{k+1} \leftarrow \text{IncrementalUpdate}_{\eta}(\theta_k)$$

$$\pi_{k+1} = \pi_{\theta_{k+1}}$$

Another Idea: Policy Gradient

Policy-centric learning: maximize the return directly
c.f. Value-centric learning: minimize the Bellman error

For $k = 1, 2, \dots$

$$\theta_{k+1} \leftarrow \theta_k + \eta \nabla_{\theta} V^{\pi_{\theta}}(\rho) \Big|_{\theta=\theta_k}$$

$$\underline{V^{\pi_{\theta}}(\rho) \triangleq \sum_s \rho(s) V^{\pi_{\theta}}(s)}$$

What are the differences between exponential weights and policy gradient?

Policy Gradient for Softmax Policy in Expert Problems

Assume full-information and fixed reward $R = (R(1), \dots, R(A))$

Let $\underline{\theta} = (\theta(1), \dots, \theta(A))$ and $\pi_{\theta}(a) = \frac{\exp(\theta(a))}{\sum_{b=1}^A \exp(\theta(b))}$

$\Rightarrow \nabla_{\theta} V^{\pi_{\theta}} = ?$

Exponential weight

$$\pi_{k+1}(a) = \frac{\pi_k(a) \exp(\eta R(a))}{\sum_b \pi_k(b) \exp(\eta R(b))}$$

??

$$V^{\pi_{\theta}} = \sum_a \pi_{\theta}(a) R(a)$$

$$PG: \underline{\theta}_{k+1} = \underline{\theta}_k + \eta \nabla_{\theta} V^{\pi_{\theta}} \Big|_{\theta=\theta_k}$$

$$\left(\nabla_{\theta} V^{\pi_{\theta}} \right)_i = \sum_a \frac{\partial}{\partial \theta_i} \left(\pi_{\theta}(a) \right) R(a) = \frac{\exp(\theta(i)) R(i)}{\sum_b \exp(\theta(b))} - \sum_a \frac{\exp(\theta(a)) \exp(\theta(i)) R(a)}{\left(\sum_b \exp(\theta(b)) \right)^2} \quad \checkmark$$

$$\text{when } a=i : \frac{\partial}{\partial \theta_i} \pi_{\theta}(a) = \frac{\partial}{\partial \theta(i)} \left[\frac{\exp(\theta(i))}{\sum_b \exp(\theta(b))} \right] = \frac{\exp(\theta(i)) \left(\sum_b \exp(\theta(b)) \right) - \exp(\theta(i)) \cdot \exp(\theta(i))}{\left(\sum_b \exp(\theta(b)) \right)^2}$$

$$\text{when } a \neq i : \frac{\partial}{\partial \theta_i} \pi_{\theta}(a) = \frac{\partial}{\partial \theta(i)} \left[\frac{\exp(\theta(a))}{\sum_b \exp(\theta(b))} \right] = \frac{0 - \exp(\theta(a)) \exp(\theta(i))}{\left(\sum_b \exp(\theta(b)) \right)^2}$$

$\frac{\partial}{\partial \theta_i} \pi_{\theta}(a)$

$$\begin{aligned}
 \underline{(\nabla_{\theta} V^{\pi_{\theta}})_i} &= \frac{\exp(\theta(i)) R(i)}{\sum_b \exp(\theta(b))} - \sum_a \frac{\exp(\theta(a)) \exp(\theta(i)) R(a)}{\left(\sum_b \exp(\theta(b))\right)^2} \\
 &= \frac{\exp(\theta(i))}{\sum_b \exp(\theta(b))} \left(R(i) - \sum_a \frac{\exp(\theta(a))}{\sum_b \exp(\theta(b))} R(a) \right) \\
 &= \pi_{\theta}(i) \left(R(i) - \sum_a \pi_{\theta}(a) R(a) \right)
 \end{aligned}$$

PG: $\theta_{k+1}(i) \leftarrow \theta_k(i) + \gamma \pi_{\theta_k}(i) \left(R(i) - \sum_a \pi_{\theta_k}(a) R(a) \right)$

Comparison between EW and PG over softmax policies

$$\theta = (\theta(a), \dots, \theta(A)), \quad \pi_{\theta}(a) = \frac{\exp(\theta(a))}{\sum_b \exp(\theta(b))}, \quad V^{\pi_{\theta}} = \sum_a \pi_{\theta}(a) R(a)$$

Policy Gradient over softmax policies

For $k = 1, 2, \dots$

$$\theta_{k+1}(a) \leftarrow \theta_k(a) + \eta \pi_{\theta_k}(a) A(a)$$

Exponential weights

For $k = 1, 2, \dots$

$$\theta_{k+1}(a) \leftarrow \theta_k(a) + \eta A(a)$$

Two Types of Policy Gradients

Policy Gradient over softmax policies

$$\theta_{k+1}(a) \leftarrow \theta_k(a) + \eta \pi_{\theta_k}(a) R(a)$$

||

$$\theta_{k+1} = \operatorname{argmax}_{\theta} \left\langle \theta - \theta_k, \nabla_{\theta} V^{\pi_{\theta_k}} \right\rangle - \frac{1}{2\eta} \|\theta - \theta_k\|^2$$

(Vanilla) Policy Gradient

Exponential weights

$$\theta_{k+1}(a) \leftarrow \theta_k(a) + \eta R(a)$$

||

$$\theta_{k+1} = \operatorname{argmax}_{\theta} \left\langle \pi_{\theta} - \pi_{\theta_k}, R \right\rangle - \frac{1}{\eta} \operatorname{KL}(\pi_{\theta}, \pi_{\theta_k})$$

Natural Policy Gradient

Approximating the NPG Update

$$\theta_{k+1} = \operatorname{argmax}_{\theta} \langle \pi_{\theta} - \pi_{\theta_k}, R \rangle - \frac{1}{\eta} \operatorname{KL}(\pi_{\theta}, \pi_{\theta_k})$$

When $\theta_{k+1} \approx \theta_k$ (i.e., when η is small), the following hold:

$$\langle \pi_{\theta} - \pi_{\theta_k}, R \rangle = V^{\pi_{\theta}} - V^{\pi_{\theta_k}} \approx (\theta - \theta_k)^{\top} \nabla_{\theta} V^{\pi_{\theta}} \Big|_{\theta=\theta_k}$$

$$\operatorname{KL}(\pi_{\theta}, \pi_{\theta_k}) \approx (\theta - \theta_k)^{\top} F_{\theta_k} (\theta - \theta_k)$$

where $F_{\theta_k} := \sum_a \pi_{\theta}(a) (\nabla_{\theta} \log \pi_{\theta}(a)) (\nabla_{\theta} \log \pi_{\theta}(a))^{\top} \Big|_{\theta=\theta_k}$

(Fisher information matrix)