

Approximate Policy Iteration and Policy-Based Learning Methods

Chen-Yu Wei

Approximate Policy Iteration (API)

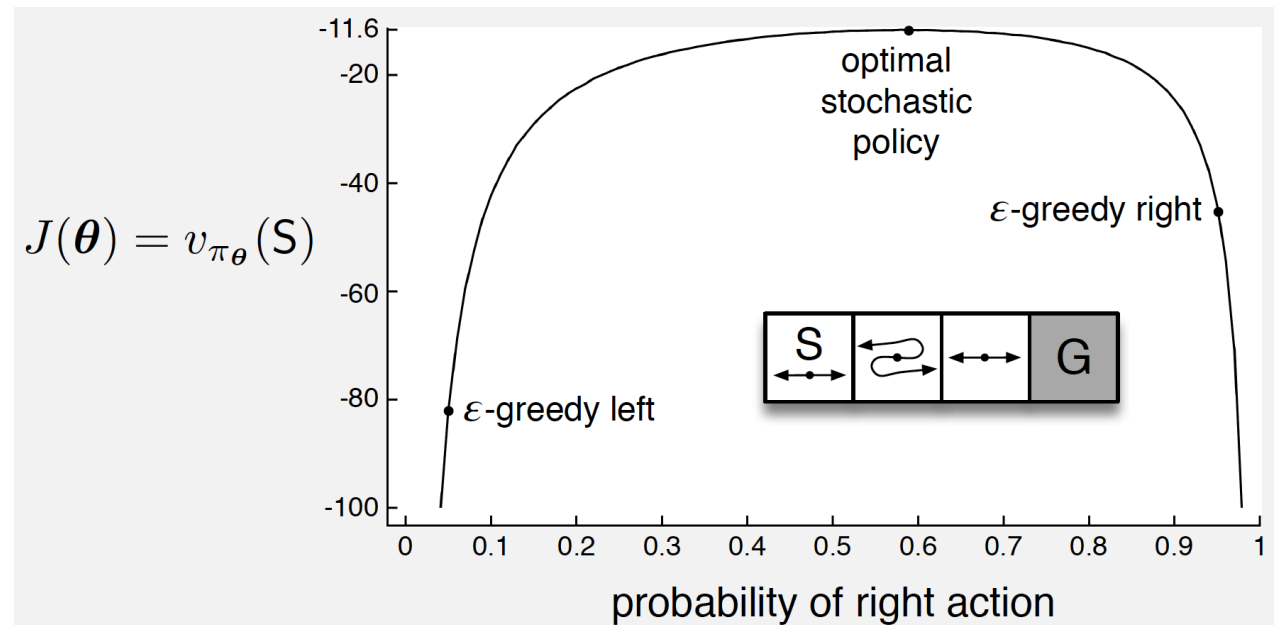
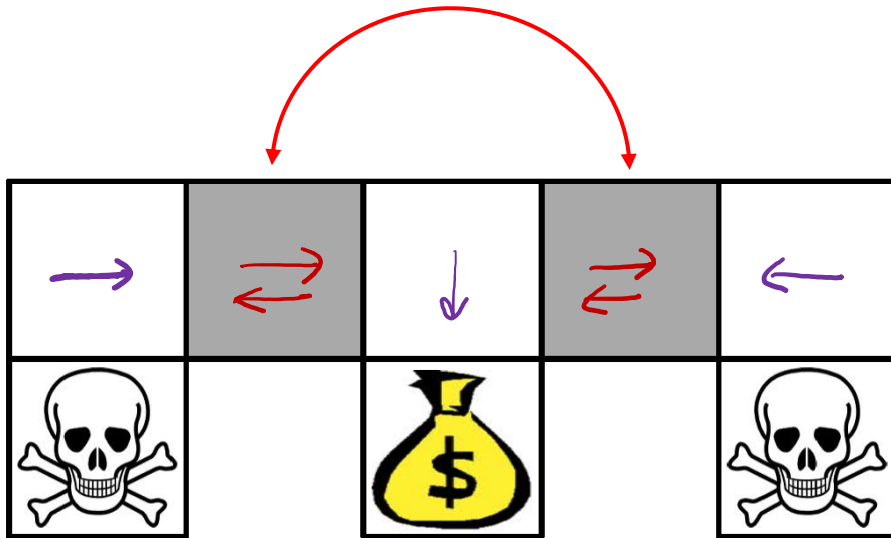
For $k = 1, 2, \dots$

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

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

Value-based : $\overset{Q^z}{Q^*}, V^z, V^* \approx \boxed{V_0}$
Policy-based : $\pi_0(a/s)$

Limitation of Value Function Approximation



Idea 1: Exponential Weights

For $k = 1, 2, \dots$

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

Perform incremental policy update such as

$$\pi_{k+1}(a|s) \propto \pi_k(a|s) \exp\left(\eta \hat{Q}_k(s, a)\right)$$

Idea 2: Policy Gradient

Parameterize policy by $\pi = \pi_\theta$

For $k = 1, 2, \dots$

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

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

V^{π_θ}

How are exponential weights and policy gradient related?

Policy Learning in the Expert Setting

Policy Gradient for Softmax Policy in Expert Problem

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

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

??

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

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

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

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

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

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

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

p6: $\theta_{k+1}(i) \leftarrow \theta_k(i) + \underbrace{\gamma \pi_{\theta_k}(i) \left(R(i) - \sum_a \pi_{\theta_k}(a) R(a) \right)}_{= A_{\theta_k}(i)}$

$$\pi_{k+1}(i) = \frac{\exp(\theta_{k+1}(i))}{\sum_b \exp(\theta_{k+1}(b))} = \frac{\underbrace{\exp(\theta_k(i))}_{\pi_k(i)} \exp(\gamma \pi_{\theta_k}(i) A_{\theta_k}(i))}{\sum_b \exp(\theta_k(b)) \exp(\gamma \pi_{\theta_k}(b) A_{\theta_k}(b))} = \frac{\pi_k(i) \exp(\gamma \pi_{\theta_k}(i) A_{\theta_k}(i))}{\sum_b \pi_k(b) \exp(\gamma \pi_{\theta_k}(b) A_{\theta_k}(b))}$$

Exponential weights:

$$A_{\pi_k}(i) = R(i) - \underbrace{\sum_a \pi_k(a) R(a)}_{\text{constant for } i}$$

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

$$\parallel$$

$$\frac{\pi_k(i) \exp(\eta R(i) - c)}{\sum_b \pi_k(b) \exp(\eta R(b) - c)} \quad \frac{\exp(-c)}{\exp(-c)}$$

PG over softmax

$$\pi_{k+1}(i) = \frac{\pi_k(i) \exp(\eta \underbrace{\pi_k(i)} A_{\pi_k}(i))}{\sum_b \pi_k(b) \exp(\eta \underbrace{\pi_k(b)} A_{\pi_k}(b))}$$

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 $k = 1, 2, \dots$

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

Exponential weights

For $k = 1, 2, \dots$

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

Experiments

Reward = [Ber(0.6), Ber(0.4)]

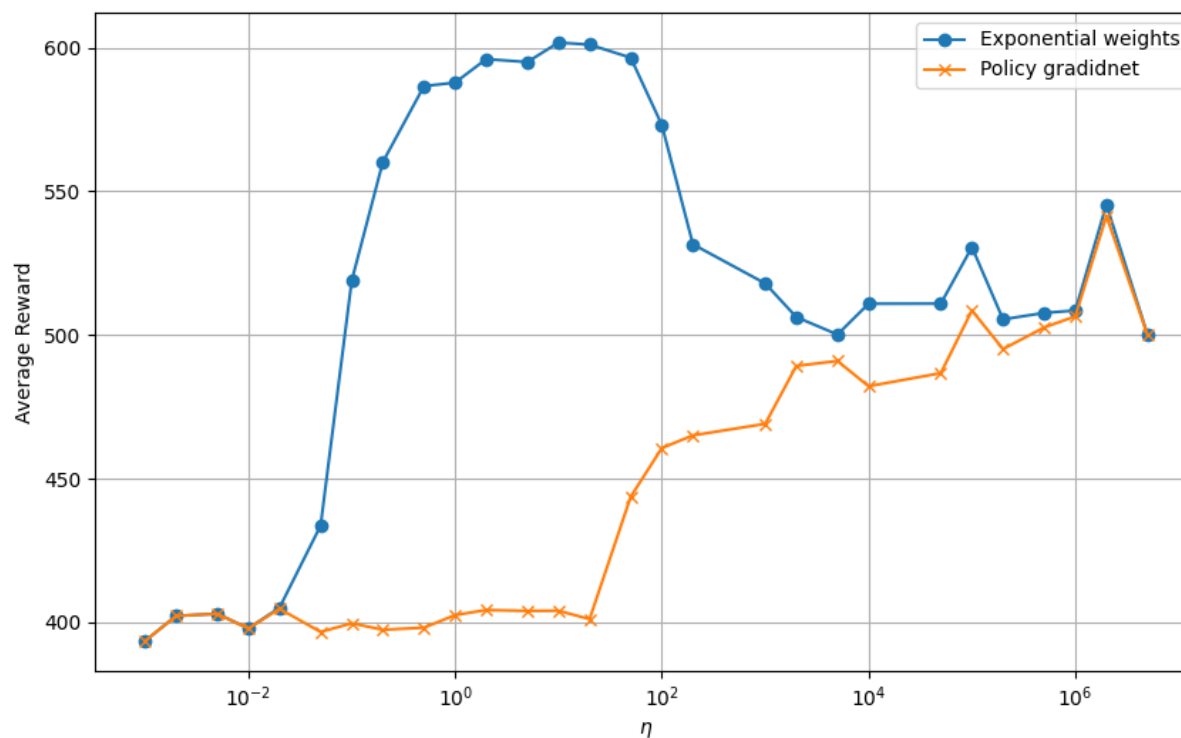
Initial policy $\pi = [0.0001, 0.9999]$

Plot total reward in 1000 rounds

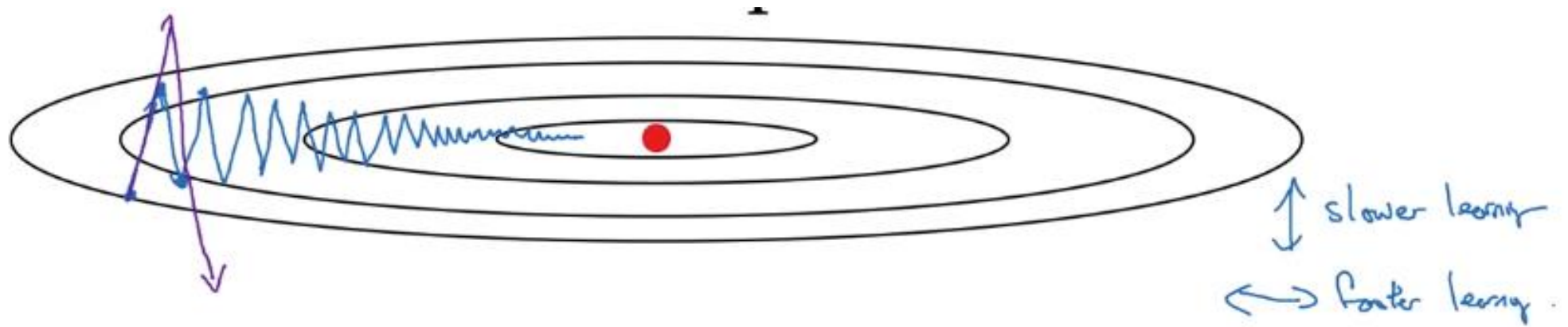
EW: $\theta_{k+1}(a) \leftarrow \theta_k(a) + \eta A_{\theta_k}(a)$

PG: $\theta_{k+1}(a) \leftarrow \theta_k(a) + \eta \pi_{\theta_k}(a) A_{\theta_k}(a)$

small eta: too slow on action 1
larger eta: too fast on action 2

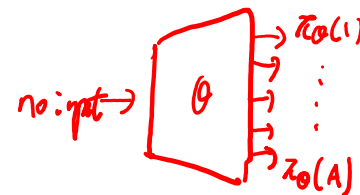


Optimization over ill-conditioned loss



<https://math.stackexchange.com/questions/2285282/relating-condition-number-of-hessian-to-the-rate-of-convergence>

Two Ideas of Policy Updates



Policy Gradient over softmax policies

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



$$\left. \nabla_{\theta} V^{\pi_{\theta}} \right|_{\theta=\theta_k} = \nabla_{\theta} V^{\pi_{\theta_k}}$$

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

Exponential weights

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



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

$$\checkmark \quad \theta_{k+1} \leftarrow \theta_k + \eta g_k$$

$$\Leftrightarrow \underset{\theta}{\operatorname{argmax}} \left\{ \langle \theta, g_k \rangle - \frac{1}{2\eta} \|\theta - \theta_k\|^2 \right\}$$

$$\checkmark \quad = \underset{\theta}{\operatorname{argmax}} \left\{ \langle \theta - \theta_k, g_k \rangle - \frac{1}{2\eta} \|\theta - \theta_k\|^2 \right\}$$

$$\Leftrightarrow \underset{\theta}{\operatorname{argmax}} \left\langle \pi_{\theta} - \pi_{\theta_k}, A_{\theta_k} \right\rangle - \frac{1}{\eta} \operatorname{KL}(\pi_{\theta}, \pi_{\theta_k})$$

$$\sum_a (\pi_{\theta}(a) - \pi_{\theta_k}(a)) \text{cost} = 0$$

$$R(a) = R(a) - \text{const}$$

Two Ideas for Function Approximation over Policies

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

(Vanilla) Policy Gradient

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

Natural Policy Gradient

Approximating the NPG Update

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

$$\begin{aligned} V^{\pi} &= \langle \pi, R \rangle \\ &= \sum_a \pi(a) R(a) \end{aligned}$$

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

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

$$\operatorname{KL}(\pi_{\theta}, \pi_{\theta_k}) \approx (\theta - \theta_k)^{\top} F_{\theta_k} (\theta - \theta_k) = \|\theta - \theta_k\|_{F_{\theta_k}}^2$$

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

(Fisher information matrix)

$$KL(\pi_\theta, \pi_{\theta+\Delta\theta}) \underset{\Delta\theta \rightarrow 0}{\approx} \frac{1}{2}(\Delta\theta)^T F_\theta (\Delta\theta) \quad \text{where } F_\theta = \sum_a \pi_\theta(a) (\nabla_\theta \log \pi_\theta(a)) (\nabla_\theta \log \pi_\theta(a))^T$$

$$KL(\pi_\theta, \pi_{\theta+\Delta\theta}) = \sum_a \pi_\theta(a) \ln \frac{\pi_\theta(a)}{\pi_{\theta+\Delta\theta}(a)}$$

$$f(\theta+\Delta\theta) \approx f(\theta) + (\nabla_\theta f(\theta))^T \Delta\theta + \frac{1}{2}(\Delta\theta)^T \underbrace{\nabla_\theta^2 f(\theta)}_{\text{Hessian}} (\Delta\theta)$$

$$= \sum_a \pi_\theta(a) \ln(\pi_\theta(a)) - \sum_a \pi_\theta(a) \ln(\pi_{\theta+\Delta\theta}(a))$$

$$\approx \sum_a \pi_\theta(a) \cancel{\ln(\pi_\theta(a))} - \sum_a \pi_\theta(a) \left(\cancel{\ln \pi_\theta(a)} + \nabla_\theta (\ln \pi_\theta(a))^T \Delta\theta + \frac{1}{2}(\Delta\theta)^T \left(\nabla_\theta^2 \ln \pi_\theta(a) \right) \Delta\theta \right)$$

$$\boxed{\nabla_\theta (\ln \pi_\theta(a)) = \frac{\nabla \pi_\theta(a)}{\pi_\theta(a)}} = \underbrace{- \sum_a \pi_\theta(a) \cdot \frac{\nabla \pi_\theta(a)^T \Delta\theta}{\pi_\theta(a)}}_{\downarrow} - \sum_a \pi_\theta(a) \cdot \frac{1}{2}(\Delta\theta)^T \left(\frac{\nabla^2 \pi_\theta(a) \pi_\theta(a) - (\nabla \pi_\theta(a))(\nabla \pi_\theta(a))^T}{(\pi_\theta(a))^2} \right) \Delta\theta$$

$$\nabla_\theta^2 (\ln \pi_\theta(a)) = \frac{(\nabla^2 \pi_\theta(a)) \pi_\theta(a) - (\nabla \pi_\theta(a))(\nabla \pi_\theta(a))^T}{(\pi_\theta(a))^2}$$

$$\begin{aligned} & - \sum_a \nabla \pi_\theta(a)^T \Delta\theta \\ &= - \nabla \left(\sum_a \pi_\theta(a) \right)^T \Delta\theta \\ &= 0 \end{aligned}$$

$$\frac{1}{2}(\Delta\theta)^T \left(\nabla^2 \sum_a \pi_\theta(a) \right) \Delta\theta$$

$$\text{For any } \theta, \sum_a \pi_\theta(a) = 1$$

NPG Updates

$$\frac{1}{2} \sum_a \pi_\theta(a) (\delta \theta)^T \left(\frac{(\nabla_\theta \pi_\theta(a)) (\nabla_\theta \pi_\theta(a))^T}{(\pi_\theta(a))^2} \right) \delta \theta \approx \frac{1}{2} (\delta \theta)^T F_\theta (\delta \theta)$$

$$= (\nabla_\theta \log \pi_\theta(a)) (\nabla_\theta \log \pi_\theta(a))^T$$

$$\theta_{k+1} = \theta_k + \eta F_{\theta_k}^{-1} \left(\nabla_\theta V^{\pi_\theta} \Big|_{\theta=\theta_k} \right)$$

$$\nabla_\theta (\log \pi_\theta(a)) = \frac{\nabla_\theta \pi_\theta(a)}{\pi_\theta(a)}$$

cf. vanilla PG: $\theta_{k+1} = \theta_k + \eta \left(\nabla_\theta V^{\pi_\theta} \Big|_{\theta=\theta_k} \right)$

NPG: $\theta_{k+1} = \underset{\theta}{\operatorname{argmax}} \left\{ \sum_a (\pi_\theta(a) - \pi_{\theta_k}(a)) R(a) - \frac{1}{\gamma} \text{KL}(\pi_\theta, \pi_{\theta_k}) \right\}$

$\approx \underset{\theta}{\operatorname{argmax}} \left\{ \langle \theta - \theta_k, \nabla_\theta V^{\pi_{\theta_k}} \rangle - \frac{1}{2\gamma} (\theta - \theta_k)^T F_{\theta_k} (\theta - \theta_k) \right\} \rightarrow W(\theta)$

$$\nabla_\theta W(\theta) = \nabla_\theta V^{\pi_{\theta_k}} - \frac{1}{\gamma} F_{\theta_k} (\theta - \theta_k) = 0 \Rightarrow \theta_{k+1} = \theta_k + \gamma F_{\theta_k}^{-1} (\nabla_\theta V^{\pi_{\theta_k}})$$

Summary: Policy Learning in the Expert Setting

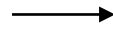
PG	NPG
$\theta_{k+1} = \operatorname{argmax}_{\theta} \langle \theta - \theta_k, \nabla_{\theta} V^{\pi_{\theta_k}} \rangle - \frac{1}{2\eta} \ \theta - \theta_k\ ^2$	$\theta_{k+1} = \operatorname{argmax}_{\theta} \langle \pi_{\theta} - \pi_{\theta_k}, R \rangle - \frac{1}{\eta} \operatorname{KL}(\pi_{\theta}, \pi_{\theta_k})$
$\theta_{k+1} = \theta_k + \eta \nabla_{\theta} V^{\pi_{\theta_k}}$	$\theta_{k+1} = \theta_k + \eta F_{\theta_k}^{-1} \nabla_{\theta} V^{\pi_{\theta_k}}$ <p>where $F_{\theta} = \mathbb{E}_{a \sim \pi_{\theta}} [(\nabla_{\theta} \log \pi_{\theta}(a))(\nabla_{\theta} \log \pi_{\theta}(a))^{\top}]$</p>
$\theta_{k+1}(a) = \theta_k(a) + \eta \pi_{\theta_k}(a) A_{\theta_k}(a)$ <p>(under direct softmax parameterization)</p>	$\theta_{k+1}(a) = \theta_k(a) + \eta A_{\theta_k}(a)$ <p>(under direct softmax parameterization)</p>

Policy Learning with Bandit Feedback

The design of EXP3

Full-information

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



Bandit

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

Inverse propensity weighting

$$\hat{r}_k(a) = \frac{r_k(a) \mathbb{I}\{a_k = a\}}{\pi_k(a)}$$

$$\hat{r}_k(a) = \frac{(r_k(a) - b - c(a)) \mathbb{I}\{a_k = a\}}{\pi_k(a)} + c(a)$$

NPG (regularization form) + Bandit Feedback

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

Use π_{θ_k} to draw $a_{k1}, a_{k2}, \dots, a_{kn}$, and get rewards $r_{k1}, r_{k2}, \dots, r_{kn}$

Approximate
$$R(a) \approx \sum_{i=1}^n \frac{(r_{ki} - b) \mathbb{I}\{a_{ki} = a\}}{\pi_{\theta_k}(a_{ki})} \quad (n = 1 \text{ recovers EXP3})$$

NPG (regularization form) + Bandit Feedback

For $k = 1, 2, \dots$

Use π_{θ_k} to draw $a_{k1}, a_{k2}, \dots, a_{kn}$, and get rewards $r_{k1}, r_{k2}, \dots, r_{kn}$

$$\text{Let } \hat{R}_k(a) = \frac{1}{n} \sum_{i=1}^n \frac{(r_{ki} - b) \mathbb{I}\{a_{ki} = a\}}{\pi_{\theta_k}(a_{ki})}$$

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

NPG (regularization form) + Bandit Feedback

For $k = 1, 2, \dots$

Use π_{θ_k} to draw $a_{k1}, a_{k2}, \dots, a_{kn}$, and get rewards $r_{k1}, r_{k2}, \dots, r_{kn}$

$$\text{Let } \hat{R}_k(a) = \frac{1}{n} \sum_{i=1}^n \frac{(r_{ki} - b) \mathbb{I}\{a_{ki} = a\}}{\pi_{\theta_k}(a_{ki})}$$

$$\theta \leftarrow \theta_k$$

Repeat m times:

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} \left(\langle \pi_{\theta} - \pi_{\theta_k}, \hat{R}_k \rangle - \frac{1}{\eta} \text{KL}(\pi_{\theta}, \pi_{\theta_k}) \right)$$

$$\theta_{k+1} \leftarrow \theta$$

PG / NPG (Gradient-Update Form) + Bandit Feedback

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

PG

$$\theta_{k+1} = \theta_k + \eta F_{\theta_k}^{-1} \left(\nabla_{\theta} V^{\pi_{\theta}} \Big|_{\theta=\theta_k} \right)$$

NPG

PG + Bandit Feedback

For $k = 1, 2, \dots$

Use π_{θ_k} to draw $a_{k1}, a_{k2}, \dots, a_{kn}$, and get rewards $r_{k1}, r_{k2}, \dots, r_{kn}$

$$\text{Let } g_k = \frac{1}{n} \sum_{i=1}^n \frac{r_{ki} - b}{\pi_{\theta_k}(a_{ki})} \left(\nabla_{\theta} \pi_{\theta}(a_{ki}) \Big|_{\theta=\theta_k} \right)$$

$$\theta_{k+1} = \theta_k + \eta g_k$$

PG + Bandit Feedback

For $k = 1, 2, \dots$

Use π_{θ_k} to draw $a_{k1}, a_{k2}, \dots, a_{kn}$, and get rewards $r_{k1}, r_{k2}, \dots, r_{kn}$

$$\text{Let } g_k = \frac{1}{n} \sum_{i=1}^n (r_{ki} - b) \nabla_{\theta} \log \pi_{\theta}(a_{ki}) \Big|_{\theta=\theta_k}$$

$$\theta_{k+1} = \theta_k + \eta g_k$$

NPG (Gradient-Update Form) + Bandit Feedback

For $k = 1, 2, \dots$

Use π_{θ_k} to draw $a_{k1}, a_{k2}, \dots, a_{kn}$, and get rewards $r_{k1}, r_{k2}, \dots, r_{kn}$

$$\text{Let } g_k = \frac{1}{n} \sum_{i=1}^n (r_{ki} - b) \nabla_{\theta} \log \pi_{\theta}(a_{ki}) \Big|_{\theta = \theta_k}$$

$$\theta_{k+1} = \theta_k + \eta F_{\theta_k}^{-1} g_k$$

Summary: Policy Learning in Bandits

PG	NPG
$\theta_{k+1} = \operatorname{argmax}_{\theta} \langle \theta - \theta_k, \nabla_{\theta} V^{\pi_{\theta_k}} \rangle - \frac{1}{2\eta} \ \theta - \theta_k\ ^2$	$\theta_{k+1} = \operatorname{argmax}_{\theta} \langle \pi_{\theta} - \pi_{\theta_k}, R \rangle - \frac{1}{\eta} \operatorname{KL}(\pi_{\theta}, \pi_{\theta_k})$
$\theta_{k+1} = \theta_k + \eta \nabla_{\theta} V^{\pi_{\theta_k}}$	$\theta_{k+1} = \theta_k + \eta F_{\theta_k}^{-1} \nabla_{\theta} V^{\pi_{\theta_k}}$ where $F_{\theta} = \mathbb{E}_{a \sim \pi_{\theta}} [(\nabla_{\theta} \log \pi_{\theta}(a))(\nabla_{\theta} \log \pi_{\theta}(a))^{\top}]$

$$\begin{aligned} \nabla_{\theta} V^{\pi_{\theta_k}} &\approx \frac{1}{n} \sum_{i=1}^n \frac{r_{ki} - b}{\pi_{\theta_k}(a_{ki})} \nabla_{\theta} \pi_{\theta}(a_{ki}) \Big|_{\theta=\theta_k} \\ &= \frac{1}{n} \sum_{i=1}^n (r_{ki} - b) \nabla_{\theta} \log \pi_{\theta}(a_{ki}) \Big|_{\theta=\theta_k} \end{aligned}$$

$$R(a) \approx \frac{1}{n} \sum_{i=1}^n \frac{(r_{ki} - b) \mathbb{I}\{a_{ki} = a\}}{\pi_{\theta_k}(a_{ki})}$$

Policy Learning in MDPs

(Full-Information Case)

Exponential Weights

For $k = 1, 2, \dots$

Perform individual exponential weight update **on all state s** :

$$\pi_{k+1}(a|s) = \frac{\pi_k(a|s) \exp(\eta Q^{\pi_k}(s, a))}{\sum_{a'} \pi_k(a'|s) \exp(\eta Q^{\pi_k}(s, a'))}$$

Analysis for Exponential Weights

Theorem.

The exponential weight algorithm guarantees

$$\sum_{k=1}^K (V^{\pi^*}(\rho) - V^{\pi_k}(\rho)) \leq \frac{1}{1-\gamma} \left(\frac{\ln A}{\eta} + \eta AK \right)$$

for any initial state distribution ρ .

Remark. It is possible to show “last-iterate convergence”

Equivalent Forms of Exponential Weights

$$\forall s, \quad \pi_{k+1}(\cdot | s) = \operatorname{argmax}_{\pi(\cdot | s)} \left\{ \underbrace{\sum_a \pi(a|s) Q^{\pi_k}(s, a)} - \frac{1}{\eta} \operatorname{KL}(\pi(\cdot | s), \pi_k(\cdot | s)) \right\}$$

$$\sum_a \pi(a|s) (Q^{\pi_k}(s, a) - b(s))$$

$$\sum_a \pi(a|s) A^{\pi_k}(s, a)$$

$$\sum_a (\pi(a|s) - \pi_k(a|s)) (Q^{\pi_k}(s, a) - b(s))$$

...

Natural Policy Gradient (Regularization Form)

$$\theta_{k+1} = \operatorname{argmax}_{\theta} \sum_s d_{\rho}^{\pi_{\theta_k}}(s) \left(\sum_a \pi_{\theta}(a|s) A^{\pi_{\theta_k}}(s, a) - \frac{1}{\eta} \operatorname{KL}(\pi_{\theta}(\cdot | s), \pi_{\theta_k}(\cdot | s)) \right)$$

Policy Gradient

$$\theta_{k+1} = \theta_k + \eta \nabla_{\theta} V^{\pi_{\theta}}(\rho) \Big|_{\theta=\theta_k} = \theta_k + \eta \sum_{s,a} d_{\rho}^{\pi_{\theta_k}}(s) \left(\nabla_{\theta} \pi_{\theta}(a|s) \Big|_{\theta=\theta_k} \right) A^{\pi_{\theta_k}}(s, a)$$

PG vs. NPG (Regularization Form)

Policy Gradient

$$\theta_{k+1} = \operatorname{argmax}_{\theta} \sum_{s,a} d_{\rho}^{\pi_{\theta_k}}(s) (\theta - \theta_k)^{\top} \left(\nabla_{\theta} \pi_{\theta_k}(a|s) \right) A^{\pi_{\theta_k}}(s, a) - \frac{1}{2\eta} \|\theta - \theta_k\|^2$$

Natural Policy Gradient

$$\theta_{k+1} = \operatorname{argmax}_{\theta} \sum_{s,a} d_{\rho}^{\pi_{\theta_k}}(s) \pi_{\theta}(a|s) A^{\pi_{\theta_k}}(s, a) - \frac{1}{\eta} \sum_s d_{\rho}^{\pi_{\theta_k}}(s) \operatorname{KL}(\pi_{\theta}(\cdot | s), \pi_{\theta_k}(\cdot | s))$$

Natural Policy Gradient (Gradient-Update Form)

$$\theta_{k+1} = \operatorname{argmax}_{\theta} \sum_{s,a} d_{\rho}^{\pi_{\theta_k}}(s) \pi_{\theta}(a|s) A^{\pi_{\theta_k}}(s, a) - \frac{1}{\eta} \sum_s d_{\rho}^{\pi_{\theta_k}}(s) \text{KL}(\pi_{\theta}(\cdot|s), \pi_{\theta_k}(\cdot|s))$$

$$\approx \operatorname{argmax}_{\theta} \sum_{s,a} d_{\rho}^{\pi_{\theta_k}}(s) (\theta - \theta_k)^{\top} \left(\nabla_{\theta} \pi_{\theta}(a|s) \Big|_{\theta=\theta_k} \right) A^{\pi_{\theta_k}}(s, a) - \frac{1}{2\eta} \sum_s d_{\rho}^{\pi_{\theta_k}}(s) (\theta - \theta_k)^{\top} F_{\theta_k}(s) (\theta - \theta_k)$$

$$= \operatorname{argmax}_{\theta} (\theta - \theta_k)^{\top} \left(\nabla_{\theta} V^{\pi_{\theta}}(\rho) \Big|_{\theta=\theta_k} \right) - \frac{1}{2\eta} (\theta - \theta_k)^{\top} F_{\theta_k} (\theta - \theta_k)$$

$$= \theta_k + \eta F_{\theta_k}^{-1} \left(\nabla_{\theta} V^{\pi_{\theta}}(\rho) \Big|_{\theta=\theta_k} \right)$$

Summary: Full-Information Policy Learning in MDPs

PG	NPG
$\operatorname{argmax}_{\theta} \langle \theta - \theta_k, \nabla_{\theta} V^{\pi_{\theta_k}}(\rho) \rangle - \frac{1}{2\eta} \ \theta - \theta_k\ ^2$	$\operatorname{argmax}_{\theta} \sum_{s,a} d_{\rho}^{\pi_{\theta_k}}(s) \pi_{\theta}(a s) A^{\pi_{\theta_k}}(s,a) - \frac{1}{\eta} \sum_s d_{\rho}^{\pi_{\theta_k}}(s) \operatorname{KL}(\pi_{\theta}(\cdot s), \pi_{\theta_k}(\cdot s))$
$\theta_k + \eta \nabla_{\theta} V^{\pi_{\theta_k}}(\rho)$	$\theta_k + \eta F_{\theta_k}^{-1} \nabla_{\theta} V^{\pi_{\theta_k}}(\rho)$ <div>Type equation here.</div> <p>where $F_{\theta} = \mathbb{E}_{s \sim d_{\rho}^{\pi_{\theta}}} \mathbb{E}_{a \sim \pi_{\theta}} [(\nabla_{\theta} \log \pi_{\theta}(a s)) (\nabla_{\theta} \log \pi_{\theta}(a s))^{\top}]$</p>
	<p>Tabular Case:</p> $\pi_{k+1}(a s) = \frac{\pi_k(a s) \exp(\eta A^{\pi_k}(s,a))}{\sum_{a'} \pi_k(a' s) \exp(\eta A^{\pi_k}(s,a'))}$

Policy Learning in MDPs

(Bandit Feedback Case)

NPG (Regularization Form)

$$\theta_{k+1} = \operatorname{argmax}_{\theta} \sum_{s,a} d_{\rho}^{\pi_{\theta_k}}(s) \pi_{\theta}(a|s) (Q^{\pi_{\theta_k}}(s,a) - b(s)) - \frac{1}{\eta} \sum_s d_{\rho}^{\pi_{\theta_k}}(s) \operatorname{KL}(\pi_{\theta}(\cdot|s), \pi_{\theta_k}(\cdot|s))$$

$$d_{\rho}^{\pi}(s) = \mathbb{E} \left[\sum_{h=1}^{\infty} \gamma^{h-1} \mathbb{I}\{s_h = s\} \mid s_1 \sim \rho, a_h \sim \pi(\cdot|s_h) \right]$$

$$Q^{\pi}(s,a) = \mathbb{E} \left[\sum_{h=1}^{\infty} \gamma^{h-1} R(s_h, a_h) \mid (s_1, a_1) = (s, a), a_h \sim \pi(\cdot|s_h) \text{ for } h \geq 2 \right]$$

NPG (Regularization Form)

$$\theta_{k+1} = \operatorname{argmax}_{\theta} \sum_{s,a} d_{\rho}^{\pi_{\theta_k}}(s) \pi_{\theta}(a|s) (Q^{\pi_{\theta_k}}(s,a) - b(s)) - \frac{1}{\eta} \sum_s d_{\rho}^{\pi_{\theta_k}}(s) \operatorname{KL}(\pi_{\theta}(\cdot|s), \pi_{\theta_k}(\cdot|s))$$

For a fixed θ , an estimator for $\sum_{s,a} d_{\rho}^{\pi_{\theta_k}}(s) \pi_{\theta}(a|s) (Q^{\pi_{\theta_k}}(s,a) - b(s))$ can be obtained as follows:

Sample a trajectory $(s_1 \sim \rho, a_1, r_1, s_2, a_2, r_2, \dots, s_{\tau}, a_{\tau}, r_{\tau})$ using policy π_{θ_k}

Define $R_h = \sum_{i=h}^{\tau} \gamma^{i-h} r_i$

Define the estimator as $\sum_{h=1}^{\tau} \gamma^{h-1} \frac{\pi_{\theta}(a_h | s_h)}{\pi_{\theta_k}(a_h | s_h)} (R_h - b(s_h))$

Similarly, $\sum_s d_{\rho}^{\pi_{\theta_k}}(s) \operatorname{KL}(\pi_{\theta}(\cdot|s), \pi_{\theta_k}(\cdot|s))$ can be estimated as $\sum_{h=1}^{\tau} \gamma^{h-1} \operatorname{KL}(\pi_{\theta}(\cdot|s_h), \pi_{\theta_k}(\cdot|s_h))$

NPG (Regularization Form)

For $k = 1, 2, \dots$

Use π_{θ_k} to collect n trajectories

$$\left(s_1^{(1)}, a_1^{(1)}, r_1^{(1)}, \dots, s_{\tau_1}^{(1)}, a_{\tau_1}^{(1)}, r_{\tau_1}^{(1)}\right), \dots, \left(s_1^{(n)}, a_1^{(n)}, r_1^{(n)}, \dots, s_{\tau_n}^{(n)}, a_{\tau_n}^{(n)}, r_{\tau_n}^{(n)}\right)$$

$$\theta_{k+1} = \operatorname{argmax}_{\theta} \left\{ \frac{1}{n} \sum_{i=1}^n \sum_{h=1}^{\tau_n} \gamma^{h-1} \frac{\pi_{\theta} \left(a_h^{(i)} \mid s_h^{(i)}\right)}{\pi_{\theta_k} \left(a_h^{(i)} \mid s_h^{(i)}\right)} \left(R_h^{(i)} - b \left(s_h^{(i)}\right)\right) \right. \\ \left. - \frac{1}{\eta} \frac{1}{n} \sum_{i=1}^n \sum_{h=1}^{\tau_n} \gamma^{h-1} \operatorname{KL} \left(\pi_{\theta} \left(\cdot \mid s_h^{(i)}\right), \pi_{\theta_k} \left(\cdot \mid s_h^{(i)}\right) \right) \right\}$$

Practical version will not include the discount factor at the front

NPG (Regularization Form)

For $k = 1, 2, \dots$

Use π_{θ_k} to collect n trajectories

$$\left(s_1^{(1)}, a_1^{(1)}, r_1^{(1)}, \dots, s_{\tau_1}^{(1)}, a_{\tau_1}^{(1)}, r_{\tau_1}^{(1)}\right), \dots, \left(s_1^{(n)}, a_1^{(n)}, r_1^{(n)}, \dots, s_{\tau_n}^{(n)}, a_{\tau_n}^{(n)}, r_{\tau_n}^{(n)}\right)$$

$$\theta_{k+1} = \operatorname{argmax}_{\theta} \left\{ \frac{1}{n} \sum_{i=1}^n \sum_{h=1}^{\tau_n} \frac{\pi_{\theta} \left(a_h^{(i)} \middle| s_h^{(i)} \right)}{\pi_{\theta_k} \left(a_h^{(i)} \middle| s_h^{(i)} \right)} \left(R_h^{(i)} - b \left(s_h^{(i)} \right) \right) \right. \\ \left. - \frac{1}{\eta} \frac{1}{n} \sum_{i=1}^n \sum_{h=1}^{\tau_n} \operatorname{KL} \left(\pi_{\theta} \left(\cdot \middle| s_h^{(i)} \right), \pi_{\theta_k} \left(\cdot \middle| s_h^{(i)} \right) \right) \right\}$$

NPG (Regularization Form)

For $k = 1, 2, \dots$

Use π_{θ_k} to collect n trajectories

$$\left(s_1^{(1)}, a_1^{(1)}, r_1^{(1)}, \dots, s_{\tau_1}^{(1)}, a_{\tau_1}^{(1)}, r_{\tau_1}^{(1)}\right), \dots, \left(s_1^{(n)}, a_1^{(n)}, r_1^{(n)}, \dots, s_{\tau_n}^{(n)}, a_{\tau_n}^{(n)}, r_{\tau_n}^{(n)}\right)$$

$$\text{Let } W_k(\theta) := \frac{1}{n} \sum_{i=1}^n \sum_{h=1}^{\tau_n} \frac{\pi_{\theta} \left(a_h^{(i)} | s_h^{(i)}\right)}{\pi_{\theta_k} \left(a_h^{(i)} | s_h^{(i)}\right)} \left(R_h^{(i)} - b \left(s_h^{(i)}\right)\right) - \frac{1}{\eta} \frac{1}{n} \sum_{i=1}^n \sum_{h=1}^{\tau_n} \text{KL} \left(\pi_{\theta} \left(\cdot | s_h^{(i)}\right), \pi_{\theta_k} \left(\cdot | s_h^{(i)}\right)\right)$$

$$\theta \leftarrow \theta_k$$

Repeat m times:

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} W_k(\theta)$$

$$\theta_{k+1} \leftarrow \theta$$

PG

$$\theta_{k+1} = \theta_k + \eta \nabla_{\theta} V^{\pi_{\theta}}(\rho) \Big|_{\theta=\theta_k} = \theta_k + \eta \sum_{s,a} d_{\rho}^{\pi_{\theta_k}}(s) \left(\nabla_{\theta} \pi_{\theta}(a|s) \Big|_{\theta=\theta_k} \right) (Q^{\pi_{\theta_k}}(s,a) - b(s))$$