# Approximate Policy Iteration and Policy-Based Learning Methods

Chen-Yu Wei

# Approximate Policy Iteration (API)

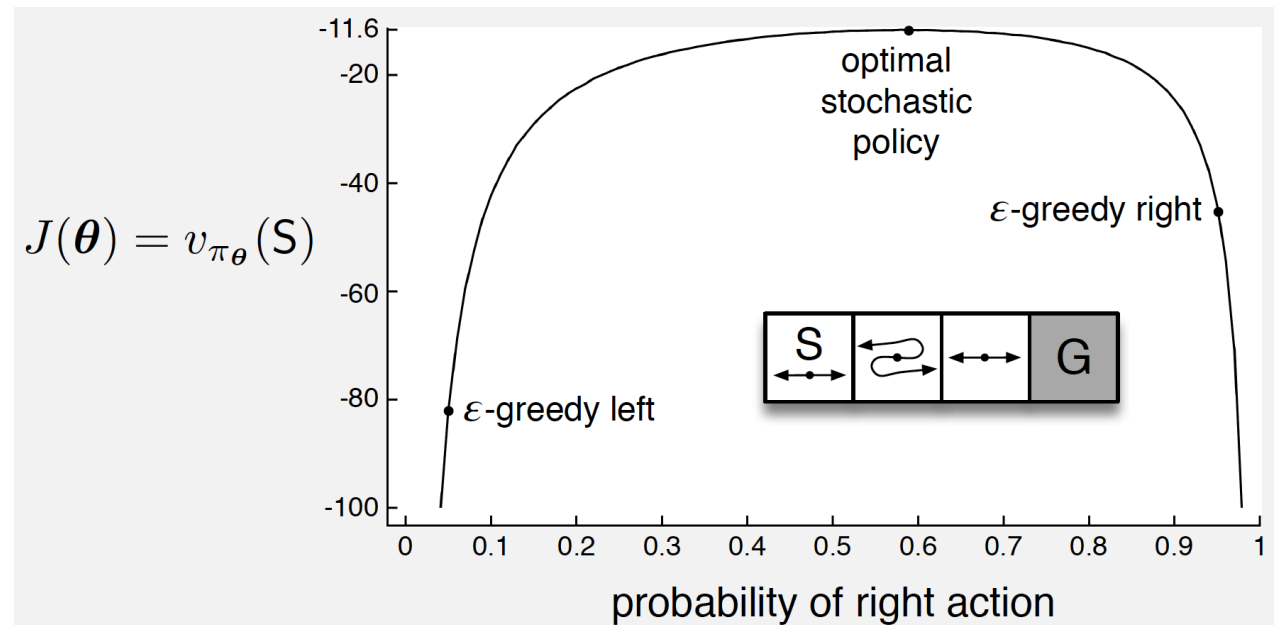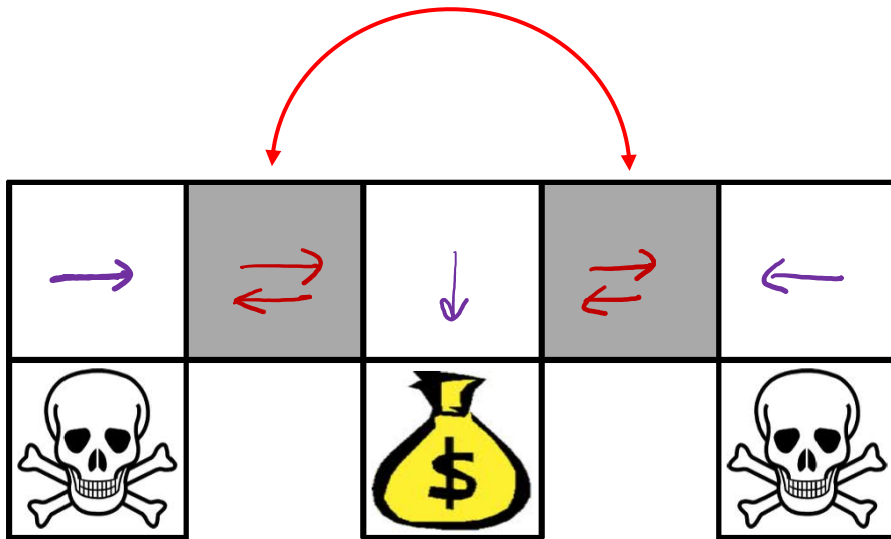For $k = 1, 2, \ldots$

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

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

$$Q^\pi$$

Value-based : $Q^*, V^\pi, V^* \approx \boxed{V_\theta}$

Policy-based : $\underline{\pi_\theta(a|s)}$

# Limitation of Value Function Approximation



$$J(\boldsymbol{\theta}) = v_{\pi_{\boldsymbol{\theta}}}(S)$$

# 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 $k = 1, \ 2, \dots$

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

$$V^{\pi_\theta}(\rho) \stackrel{\Delta}{=} \sum_s \rho(s) V^{\pi_\theta}(s)$$

$$V^{\pi_\theta}$$

How are exponential weights and policy gradient related?

# Policy Gradient 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 $\theta = (\theta(1), \dots, \theta(A))$ and $\pi_\theta(a) = \dfrac{\exp(\theta(a))}{\sum_{b=1}^{A} \exp(\theta(b))}$

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

Exponential weight

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

??

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

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

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

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

when $a = i$:

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

when $a \neq i$:

$$\frac{\partial}{\partial \theta_i} \pi_\theta(a) = \frac{\partial}{\partial \theta(i)}\left(\frac{\exp(\theta(a))}{\sum_b \exp(\theta(b))}\right) = \frac{0 - \exp(\theta(a)) \exp(\theta(i))}{\left(\sum_b \exp(\theta(b))\right)^2}$$

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

$$= \frac{\exp(\theta(i))}{\sum_b \exp(\theta(b))}\left( R(i) - \sum_a \frac{\exp(\theta(a))}{\sum_b \exp(\theta(b))} R(a)\right)$$

$$= \pi_\theta(i)\left( R(i) - \sum_a \pi_\theta(a) R(a)\right)$$

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

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

$$A_{\pi_k}(i) = R(i)\left(-\sum_a \pi_k(a)R(a)\right) \rightarrow \text{constant for } i$$

Exponential weights:

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

$$\|$$

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

PG over softmax

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

# Comparison between EW and PG over softmax policies

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

**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

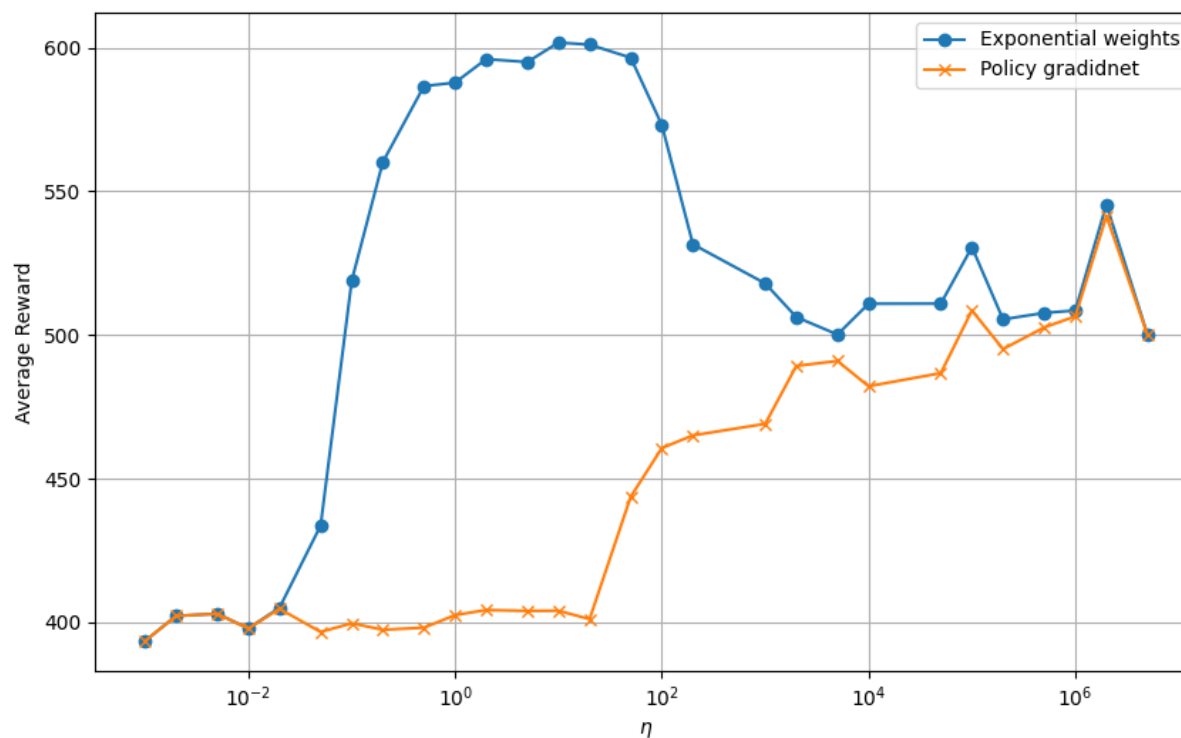Reward = $[\text{Ber}(0.6), \text{Ber}(0.4)]$

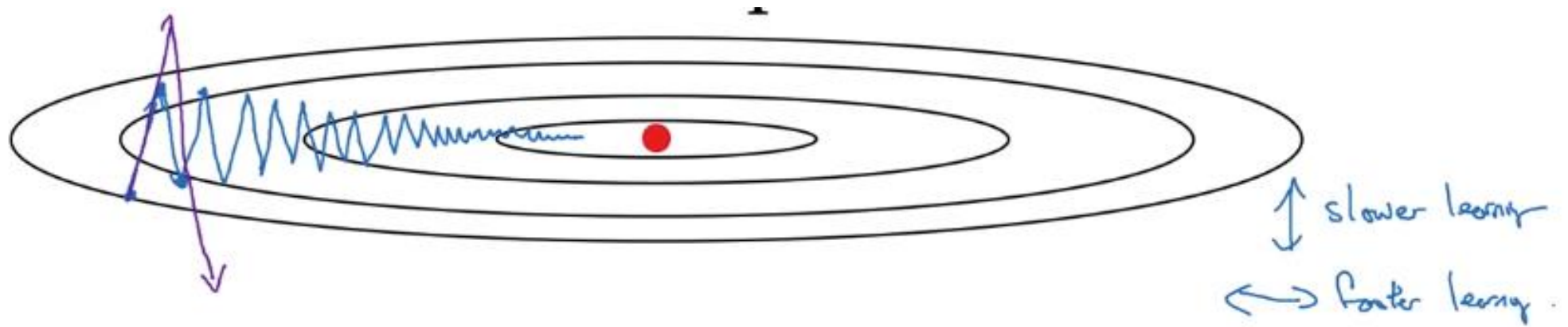Initial policy $\pi = [0.0001, 0.9999]$

Plot total reward in 1000 rounds

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

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

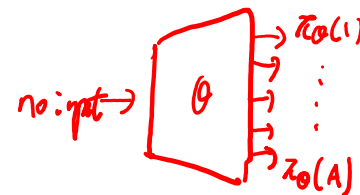small eta : too slow on action 1
larger eta : too fast on action 2

# Optimization over ill-conditioned loss

# Two Ideas of Policy Updates



## Policy Gradient over softmax policies

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

$$\nabla_\theta V^{\pi_\theta}\big|_{\theta = \theta_k} = \nabla_\theta V^{\pi_{\theta_k}}$$

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

## Exponential weights

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

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

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

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

$$(\Rightarrow) \quad \theta_{k+1} = \underset{\theta}{\operatorname{argmax}} \left\{ \langle \theta, g_k \rangle - \frac{1}{2\eta} \|\theta - \theta_k\|^2 \right\}$$

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

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

$$\sum_a \left(\pi_\theta(a) - \pi_{\theta_k}(a)\right) \text{const} = 0$$

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

# 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**

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

**Natural Policy Gradient**

# Approximating the NPG Update

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

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

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

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

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

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

**(Fisher information matrix)**

$$KL\left(\pi_\theta, \pi_{\theta+\Delta\theta}\right) \approx \frac{1}{2}(\Delta\theta)^T F_\theta (\Delta\theta) \qquad \text{where } F_\theta = \sum_a \pi_\theta(a)\left(\nabla_\theta \log \pi_\theta(a)\right)\left(\nabla_\theta \log \pi_\theta(a)\right)^T$$

$$\underset{\Delta\theta \to 0}{\downarrow}$$

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

$$f(\theta+\Delta\theta) \approx f(\theta) + \left(\nabla_\theta f(\theta)\right)^T \Delta\theta + \frac{1}{2}(\Delta\theta)^T \nabla_\theta^2 f(\theta)(\Delta\theta)$$

Hessian

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

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

$$\boxed{\nabla_\theta\left(\ln \pi_\theta(a)\right) = \frac{\nabla \pi_\theta(a)}{\pi_\theta(a)}}$$

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

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

$$\downarrow$$

$$-\sum_a \nabla \pi_\theta(a)^T \Delta\theta$$

$$= -\nabla\left(\sum_a \pi_\theta(a)\right)^T \Delta\theta$$

$$= 0$$

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

$$\underbrace{\qquad}_{1}$$

For any $\theta$. $\sum_a \pi_\theta(a) = 1$

# NPG Updates

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

$$\underbrace{}_{} = (\nabla_\theta \log \pi_\theta(a))(\nabla_\theta \log \pi_\theta(a))^T$$

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

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

*cf.* vanilla PG:

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

NPG:

$$\theta_{k+1} = \underset{\theta}{\arg\max} \left\{ \sum_a \left( \pi_\theta(a) - \pi_{\theta_k}(a) \right) R(a) - \frac{1}{\eta} KL\left( \pi_\theta, \pi_{\theta_k} \right) \right\}$$

$$\approx \underset{\theta}{\arg\max} \left\{ \langle \theta - \theta_k, \nabla_\theta V^{\pi_{\theta_k}} \rangle - \frac{1}{2\eta}(\theta - \theta_k)^T F_{\theta_k} (\theta - \theta_k) \right\} \longrightarrow W(\theta)$$

$$\nabla_\theta W(\theta) = \nabla_\theta V^{\pi_{\theta_k}} - \frac{1}{\eta} F_{\theta_k}(\theta - \theta_k) = 0 \Rightarrow \theta_{k+1} = \theta_k + \eta F_{\theta_k}^{-1}\left( \nabla_\theta V^{\pi_{\theta_k}} \right)$$
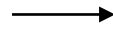
# Summary: Policy Learning in the Expert Setting

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

# Policy Learning with Bandit Feedback

# The design of EXP3

**Full-information**

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

$\longrightarrow$

**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

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

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

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

# NPG (regularization form) + Bandit Feedback

For $k = 1, \ 2, \dots$

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

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

$$\theta_{k+1} = \underset{\theta}{\arg\max} \ \langle \pi_\theta - \pi_{\theta_k}, \hat{R}_k \rangle - \frac{1}{\eta} \mathrm{KL}\big(\pi_\theta, \pi_{\theta_k}\big)$$

# NPG (regularization form) + Bandit Feedback

For $k = 1, 2, \ldots$

Use $\pi_{\theta_k}$ to draw $a_{k1}, a_{k2}, \ldots, a_{kn}$, and get rewards $r_{k1}, r_{k2}, \ldots, r_{kn}$

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

$\theta \leftarrow \theta_k$

Repeat $m$ times:

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

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

# 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 $\pi_{\theta_k}$ to draw $a_{k1}, a_{k2}, \dots, a_{kn}$, and get rewards $r_{k1}, r_{k2}, \dots, r_{kn}$

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

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

# NPG (Gradient-Update Form) + Bandit Feedback

For $k = 1, \ 2, \dots$

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

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

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

# Summary: Policy Learning in Bandits

| PG | NPG |
|---|---|
| $\theta_{k+1} = \underset{\theta}{\text{argmax}} \left\langle \theta - \theta_k, \nabla_\theta V^{\pi_{\theta_k}} \right\rangle - \dfrac{1}{2\eta} \|\theta - \theta_k\|^2$ | $\theta_{k+1} = \underset{\theta}{\text{argmax}} \left\langle \pi_\theta - \pi_{\theta_k}, \boxed{R} \right\rangle - \dfrac{1}{\eta} \text{KL}(\pi_\theta, \pi_{\theta_k})$ |
| $\theta_{k+1} = \theta_k + \eta \boxed{\nabla_\theta V^{\pi_{\theta_k}}}$ | $\theta_{k+1} = \theta_k + \eta F_{\theta_k}^{-1} \nabla_\theta V^{\pi_{\theta_k}}$ <br> where $F_\theta = \mathbb{E}_{a \sim \pi_\theta}[(\nabla_\theta \log \pi_\theta(a))(\nabla_\theta \log \pi_\theta(a))^\top]$ |

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

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