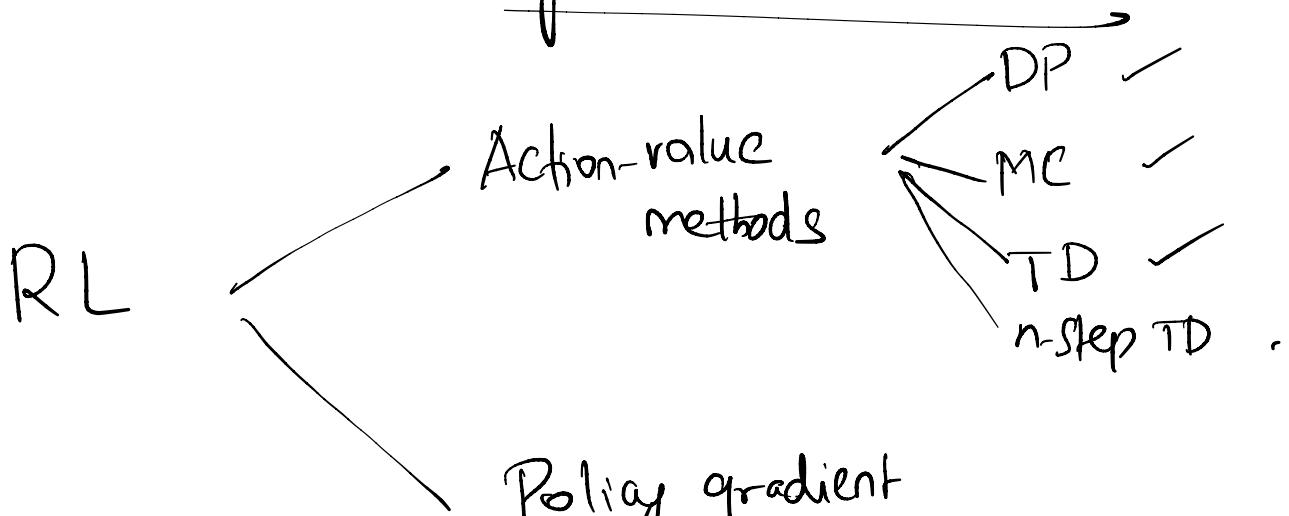


Lecture -08

Policy Gradient Methods



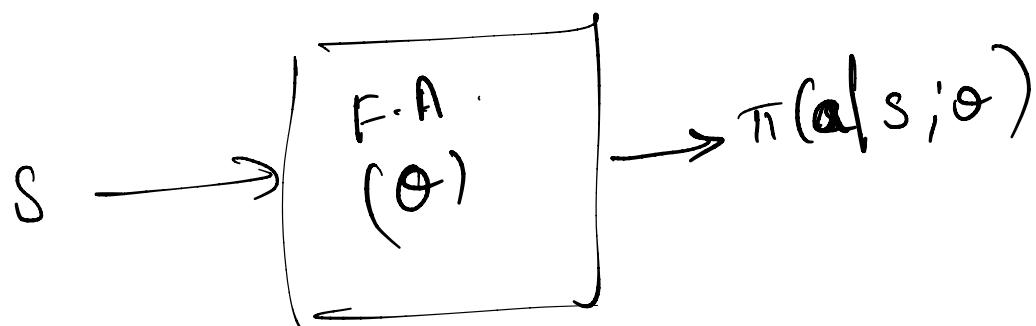
Policy gradient methods :-

Parameter of the FA
 $\hat{V}(s; w)$

$$\pi(a|s; \theta) = \Pr\{A_t = a | S_t = s, \Theta_t = \theta\}$$

θ - Parameter for the Policy approximator.

w - Param. for the value-fn approx.



performance measure : $J(\theta) = V_{\pi_\theta}(s_0)$

goal: maximize the perf.

$$\theta_{t+1} = \theta_t + \alpha \nabla J(\theta_t)$$

→ Actor-Critic methods

↑
Policy ↑
Value fn.

Advantages of PG methods:-

numerical pref $h(s, a; \theta) \in \mathbb{R}$

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

Softmax
in
action
pref.

$(0.999, 0.0001, 0.0001)$

$$h(s, a; \theta) = \theta^T x(s, a)$$

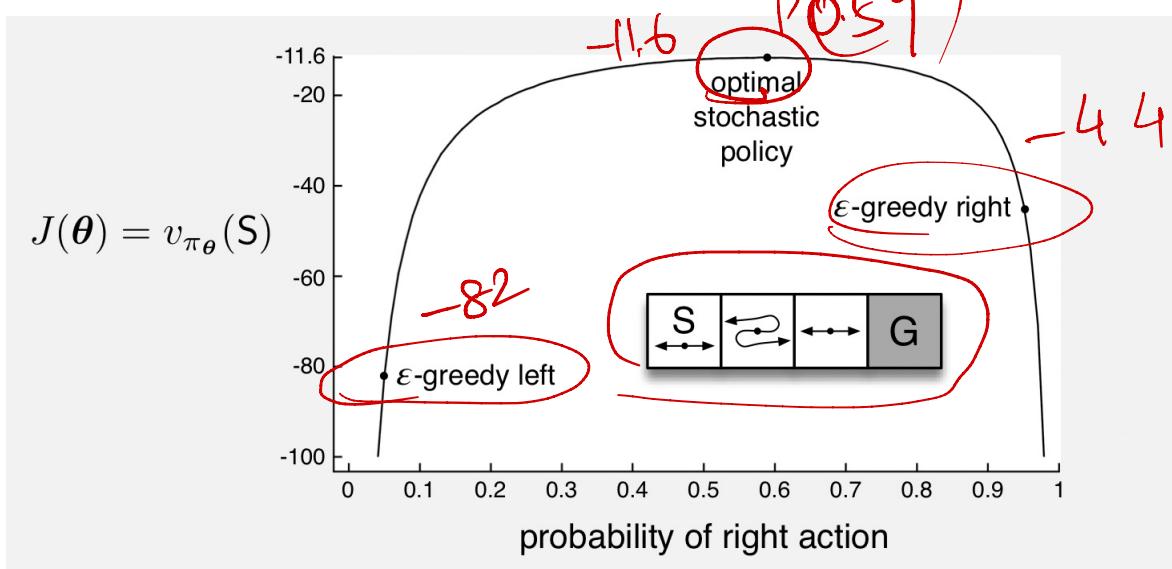
\uparrow
feature vector for s .

- ① Stochastic policy can approach deterministic policy if needed.

$$q(s, a)$$

- ② \rightarrow enable the selection of actions with arbitrary policy

Reward = -1 all the time



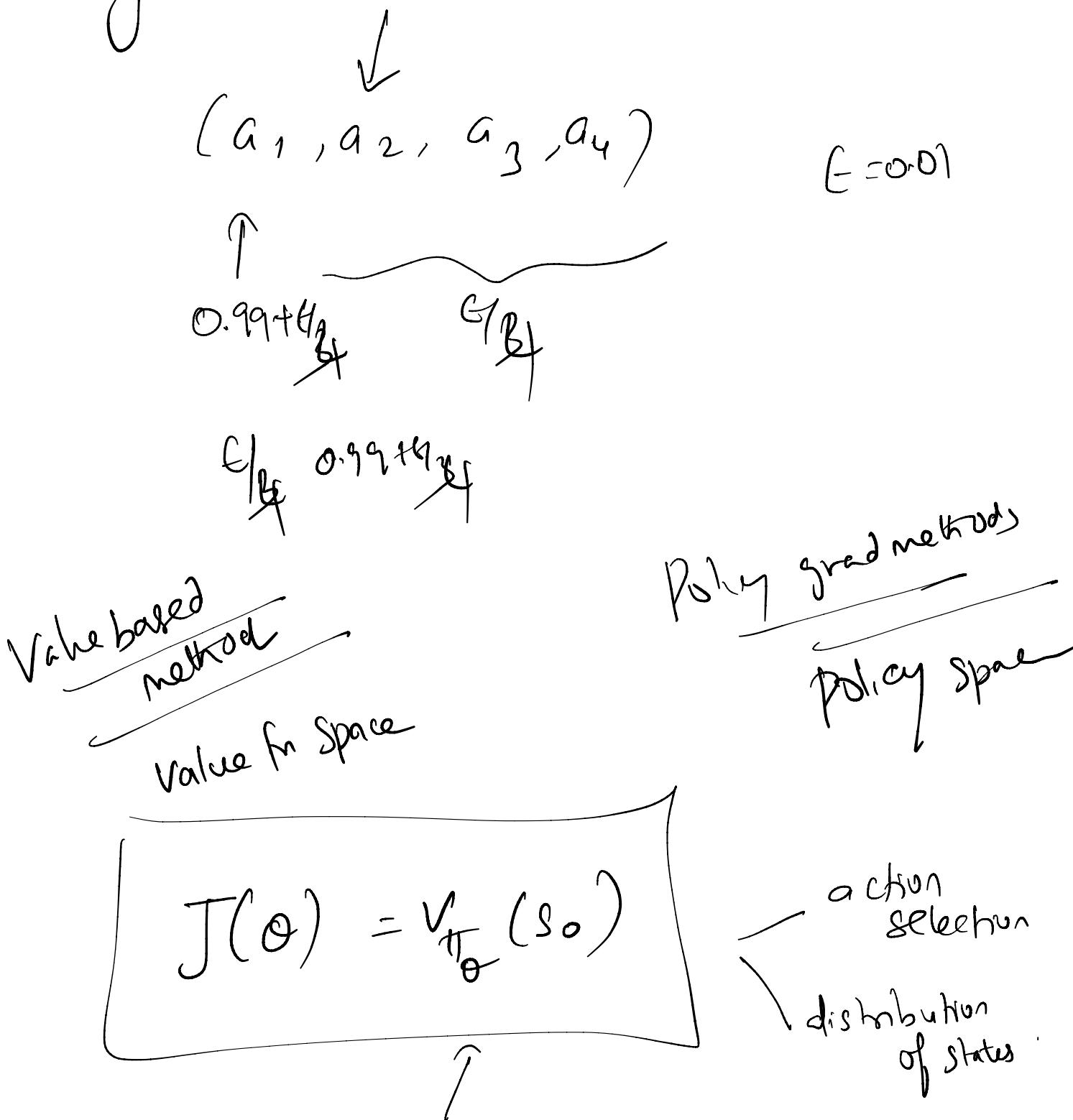
$$x(s, \text{right}) = [1, 0]^T$$

$$x(s, \text{left}) = [0, 1]^T$$

③ Policy may be simpler to approximate than value.

④ Policy Parameterization \rightarrow injecting prior knowledge.

Policy Gradient Theorem:-



$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a q_{\pi}(s, a) \nabla \pi(a|s; \theta)$$

$$\propto \mathbb{E}_{\pi} \left[\sum_a q_{\pi}(s, a) \nabla \pi(a|s; \theta) \right]$$

Proof: $\gamma = 1$

$$\begin{aligned}\nabla V_{\pi}(s) &= \nabla \left[\sum_a \pi(a|s) q_{\pi}(s, a) \right] \\ &= \sum_a \left[\nabla (\pi(a|s) q_{\pi}(s, a)) \right] \\ &= \sum_a \left[\nabla \pi(a|s) q_{\pi}(s, a) + \pi(a|s) \nabla q_{\pi}(s, a) \right]\end{aligned}$$

$$= \sum_a \left[\nabla \pi(a|s) q_{\pi}(s, a) + \pi(a|s) \right]$$
$$\nabla \sum_{s', r} P(s', r | s, a) (r + \gamma V_{\pi}(s'))$$

$$= \sum_a \left[\nabla \pi(a|s) q_{\pi}(s, a) + \pi(a|s) \right]$$
$$\sum_{s', r} P(s', r | s, a) \nabla V_{\pi}(s')$$

$$= \sum_a \left[\nabla \pi(a|s) q_{\pi}(s, a) + \pi(a|s) \right]$$
$$\sum_{s'} P(s' | s, a) \nabla V_{\pi}(s')$$

$$\nabla V^\pi(s) = \phi(s)$$

$$\sum_a \left[\nabla_{\pi}(a|s) q_\pi(s, a) + \pi(a|s) \right] \xrightarrow{\text{red lines}} \sum_{s'} p(s'|s, a) \nabla V_\pi(s')$$

$P^\pi(s \rightarrow x, k)$ Prob of reaching state
x from s , after k steps .

$$k=0, \quad P^\pi(s \rightarrow s, k=0) = 1$$

$$k=1, \quad P^\pi(s \rightarrow s', k=1) = \sum_a \pi(a|s) P(s'|s, a)$$

$s \xrightarrow{k+1} x$ after $k+1$ steps .

$s \xrightarrow{k} s' \xrightarrow{} x$

$$P^\pi(s \rightarrow x, k+1) = \sum_{s'} P^\pi(s \rightarrow s', k) \cdot P^\pi(s' \rightarrow x, 1)$$

$$\nabla V^\pi(s) = \phi(s)$$

$$\sum_a \left[\nabla_{\pi}(a|s) q_{\pi}(s,a) + \pi(a|s) \right] \equiv \sum_{s'} p(s'|s,a) \nabla V_\pi(s')$$

$$= \phi(s) + \sum_a \pi(a|s) \sum_{s'} p(s'|s,a) \nabla V_\pi(s')$$

$$= \phi(s) + \sum_{s'} \underbrace{\sum_a \pi(a|s) p(s'|s,a) \nabla V_\pi(s')}$$

$$\nabla V_\pi(s) = \phi(s) + \sum_{s'} P^\pi(s \rightarrow s', 1) \boxed{\nabla V_\pi(s')}$$

$$= \phi(s) + \sum_{s'} P^\pi(s \rightarrow s', 1) \left[\phi(s') + \sum_{s''} P^\pi(s' \rightarrow s'', 1) \nabla V_\pi(s'') \right]$$

$$= \phi(s) + \sum_{s'} P^\pi(s \rightarrow s', 1) \phi(s') + \sum_{s''} P^\pi(s \rightarrow s'', 2) \nabla V_\pi(s'') \equiv$$

$$= \phi(s) + \sum_{s'} P^{\pi}(s \rightarrow s', 1) \phi(s') + \\ \sum_{s''} P^{\pi}(s \rightarrow s'', 2) \nabla V_{\pi}(s'') \\ \underline{\underline{=}}$$

$$= \phi(s) + \sum_{s'} P^{\pi}(s \rightarrow s', 1) \phi(s') \\ + \sum_{s''} P^{\pi}(s \rightarrow s'', 2) \phi(s'') \\ + \sum_{s'''} P^{\pi}(s \rightarrow s''', 3) \nabla V_{\pi}(s''')$$

- - - - .

$$\nabla J(\theta) = \nabla V^{\pi}(s_0) = \sum_{x \in S} \sum_{k=0}^{\infty} P^{\pi}(s \rightarrow x, k) \phi(x)$$

$$= \sum_{x \in S} n(x) \phi(x) \\ \underline{\underline{=}}$$

$$= \left(\sum_S n(s) \right) \sum_S \frac{n(s)}{\sum_S n(s)} \phi(s)$$

$$\propto \sum_S \frac{n(s)}{\sum_S n(s)} \phi(s)$$

$$\nabla J(\theta) \propto \sum_s \mu(s) \phi(s)$$

$$= \sum_s \mu(s) \sum_a \nabla \pi(a|s) q_{\pi}(s, a)$$

Monte Carlo Policy Gradient :-

$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a q_{\pi}(s, a) \nabla \pi(a|s; \theta)$$

$$= \mathbb{E}_{\pi} \left[\sum_a q_{\pi}(s, a) \nabla \pi(a|s; \theta) \right]$$

①

$$\theta_{t+1} = \theta_t + \alpha \sum_a \hat{q}(s_t, a; \omega) \nabla \pi(a|s_t; \theta)$$

② REINFORCE

(Williams, 1992)

REINFORCE:

$$\begin{aligned}\nabla J(\theta) &\propto \mathbb{E}_{\pi} \left[\sum_a q_{\pi}(s_t, a) \nabla \pi(a | s_t, \theta) \right] \\&= \mathbb{E}_{\pi} \left[\sum_a \frac{\pi(a | s_t, \theta) \cdot q_{\pi}(s_t, a) \nabla \pi(a | s_t, \theta)}{\pi(a | s_t, \theta)} \right] \\&= \mathbb{E}_{\pi} \left[q_{\pi}(s_t, A_t) \frac{\nabla \pi(A_t | s_t; \theta)}{\pi(A_t | s_t; \theta)} \right] \\&= \mathbb{E}_{\pi} \left[G_t \frac{\nabla \pi(A_t | s_t; \theta)}{\pi(A_t | s_t; \theta)} \right]\end{aligned}$$

$$\theta_{t+1} = \theta_t + \alpha G_t \frac{\nabla \pi(A_t | s_t; \theta)}{\pi(A_t | s_t; \theta)}$$

$$\theta_{t+1} = \theta_t + \alpha G_t \nabla \ln \pi(A_t | s_t; \theta)$$

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

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

Algorithm parameter: step size $\alpha > 0$

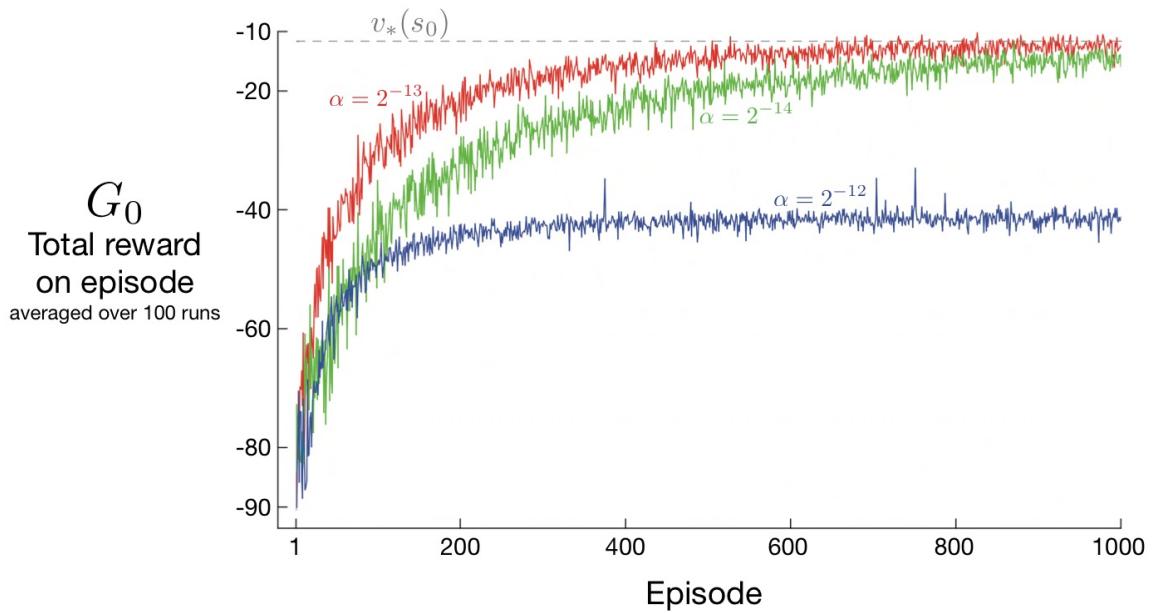
Initialize policy parameter $\theta \in \mathbb{R}^{d'}$ (e.g., to $\mathbf{0}$)

Loop forever (for each episode):

Generate an episode $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot| \cdot, \theta)$

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

$$\begin{aligned} G &\leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k \\ \theta &\leftarrow \theta + \alpha \gamma^t G \nabla \ln \pi(A_t | S_t, \theta) \end{aligned} \tag{G_t}$$



REINFORCE

$$\theta_{t+1} = \theta_t + \alpha G_t \nabla \ln \pi(A_t | s_t; \theta_t)$$

REINFORCE with Baseline:

$$\nabla J(\theta) \propto \sum_S M(s) \sum_a (q_{\pi}(s, a) - b(s)) \nabla \pi(a|s; \theta)$$

$$\begin{aligned}
 \sum_a b(s) \nabla \pi(a|s; \theta) &= b(s) \sum_a \nabla \pi(a|s; \theta) \\
 &\quad \text{---} \\
 &= b(s) \nabla 1 = 0.
 \end{aligned}$$

$$\theta_{t+1} = \theta_t + \alpha (G_t - b(s_t)) \ln \pi(A_t | s_t; \theta_t)$$

$$\theta_{t+1} = \theta_t + \alpha (G_t - b(s_t)) \ln \pi(A_t | s_t, \theta_t)$$

$$= \theta_t + \alpha (G_t - \hat{V}(s_t; \omega)) \ln \pi(A_t | s_t, \theta_t)$$

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 $\hat{v}(s, w)$

Algorithm parameters: step sizes $\alpha^\theta > 0$, $\alpha^w > 0$

Initialize policy parameter $\theta \in \mathbb{R}^{d'}$ and state-value weights $w \in \mathbb{R}^d$ (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 $\pi(\cdot| \cdot, \theta)$

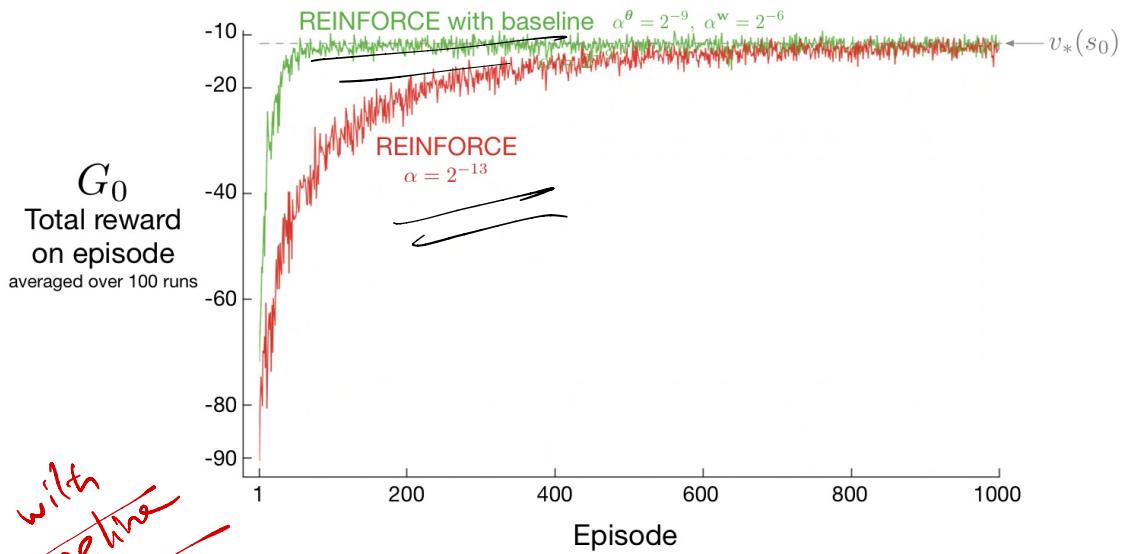
Loop for each step of the episode $t = 0, 1, \dots, T - 1$:

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

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

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

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



REINFORCE with baseline

$$\theta_{t+1} = \theta_t + \alpha (G_t - \hat{v}(s_t; \omega)) \nabla \ln \pi(a_t | s_t; \theta)$$

Actor-Critic Method:

$$\theta_{t+1} = \theta_t + \alpha (G_{[t:t+1]} - \hat{v}(s_t; \omega)) \nabla \ln \pi(a_t | s_t; \theta)$$

$$= \theta_t + \alpha \left(R_{t+1} + \gamma \hat{v}(s_{t+1}; \omega) - \hat{v}(s_t; \omega) \right) \nabla \ln \pi$$

$$\theta_{t+1} = \theta_t + \alpha \delta_t \nabla \ln \pi(a_t | s_t; \theta_t)$$

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 $\hat{v}(s, w)$

Parameters: step sizes $\alpha^\theta > 0, \alpha^w > 0$

Initialize policy parameter $\theta \in \mathbb{R}^{d'}$ and state-value weights $w \in \mathbb{R}^d$ (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

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

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

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

$I \leftarrow \gamma I$

$S \leftarrow S'$

PG for episodic problems: -

Continuing prob ~~set with setting~~ with no discounting γ

Average reward setting

$$J(\theta) = r(\pi)$$

$$= \lim_{h \rightarrow \infty} \frac{1}{h} \sum_{t=1}^h \mathbb{E}[R_t | s_0, a_0: t-1]$$

$$= \lim_{t \rightarrow \infty} \mathbb{E}[R_t | s_0, a_0: t-1]$$

$$= \sum_s \mu(s) \sum_a \pi(a|s) \sum_{s'} p(s'|s, a) r$$

where μ — steady state distrib.

$$\mu(s) = \lim_{t \rightarrow \infty} \Pr[S_t = s | A_0: t-1]$$

$$\sum_s \mu(s) \sum_a \pi(a|s, \theta) p(s'|s, a) = \mu(s')$$

$$G_t = R_{t+1} - r(\pi) + R_{t+2} - r(\pi) + R_{t+3} - r(\pi) + \dots$$

~~P GT (Continuing case)~~

$$J(\theta) = \delta(\pi)$$

$$\nabla V_\pi(s) = \nabla \left[\sum_a \pi(a|s) q_\pi(s, a) \right]$$

$$= \sum_a \left[\nabla \pi(a|s) q_\pi(s, a) + \pi(a|s) \nabla q_\pi(s, a) \right]$$

$$= \sum_a \left[\nabla \pi(a|s) q_\pi(s, a) + \pi(a|s) \cdot \right.$$

$$\left. \nabla \sum_{s', r} p(s', r|s, a) (\gamma - r(\theta) + V_\pi(s')) \right]$$

$$\nabla V_\pi(s) \equiv \sum_a \left[\nabla \pi(a|s) q_\pi(s, a) + \pi(a|s) \cdot \right. \\ \left. - \nabla r(\theta) + \sum_{s'} p(s'|s, a) \nabla V_\pi(s') \right]$$

$$\nabla r(\theta) = \sum_a \left[\nabla \pi(a|s) q_\pi(s, a) + \pi(a|s) \sum_{s'} p(s'|s, a) \nabla V_\pi(s') \right]$$

$$- \nabla V_\pi(s)$$

$$\begin{aligned}
\nabla r(\theta) &= \\
&\sum_s \mu(s) \sum_a \left[\nabla \pi(a|s) q_\pi(s,a) + \pi(a|s) \sum_{s'} p(s'|s,a) \nabla v_\pi(s') \right] \\
&\quad - \nabla v_\pi(s) \\
&= \sum_s \mu(s) \sum_a \nabla \pi(a|s) q_\pi(s,a) + \\
&\quad \sum_s \mu(s) \sum_a \pi(a|s) \sum_{s'} p(s'|s,a) \nabla v_\pi(s') \\
&\quad - \sum_s \mu(s) \nabla v_\pi(s) \\
&= \sum_s \mu(s) \sum_a \nabla \pi(a|s) q_\pi(s,a) + \\
&\quad \underbrace{\sum_{s'} \sum_s \mu(s) \sum_a \pi(a|s) p(s'|s,a) \nabla v_\pi(s')}_{-\sum_s \mu(s) \nabla v_\pi(s)} \\
&= \sum_s \mu(s) \sum_a \nabla \pi(a|s) q_\pi(s,a) + \\
&\quad \cancel{\sum_{s'} \mu(s') \nabla v_\pi(s')} \\
&\quad - \cancel{\sum_s \mu(s) \nabla v_\pi(s')} \\
\nabla J(\theta) &= \sum_s \mu(s) \sum_a \nabla \pi(a|s) q_\pi(s,a) //
\end{aligned}$$

on-policy PG algorithms

Off-Policy PG algorithms?

π - target policy
 b - behaviour policy.

$$J(\theta) = \sum_s \mu^b(s) \sum_a q_{\pi}(s, a) \pi(a|s; \theta)$$

$$= \mathbb{E}_{\substack{s \sim \mu^b(s)}} \left[\sum_a \pi(a|s; \theta) q_{\pi}(s, a) \right]$$

↑

Off-Policy PG : (Degris, White & Sutton, 2012)

$$\nabla J(\theta) = \nabla \mathbb{E}_{s \sim \mu^b} \left[\sum_a \pi(a|s; \theta) q_{\pi}(s, a) \right]$$

$$= \mathbb{E}_{\substack{s \sim \mu^b}} \left[\sum_a q_{\pi}(s, a) \nabla \pi(a|s; \theta) + \pi(a|s; \theta) \nabla q_{\pi}(s, a) \right]$$

$$\approx \mathbb{E}_{s \sim \mu^b} \left[\sum_a q_{\pi}(s, a) \nabla \pi(a|s; \theta) \right]$$

$$= \mathbb{E}_{\substack{s \sim \mu^b}} \left[\sum_a b(a|s) \frac{\pi(a|s)}{b(a|s)} q_{\pi}(s, a) \frac{\nabla \pi(a|s)}{\pi(a|s)} \right]$$

$$= \mathbb{E}_b \left[\frac{\pi(a|s; \phi)}{b(a|s)} q_{\pi}(s, a) \nabla \ln \pi(a|s) \right]$$

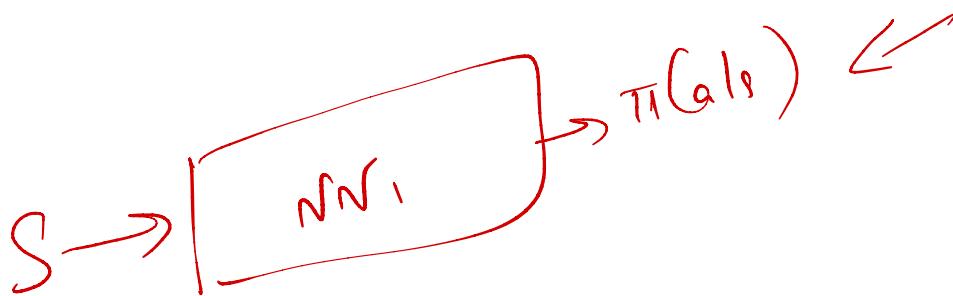
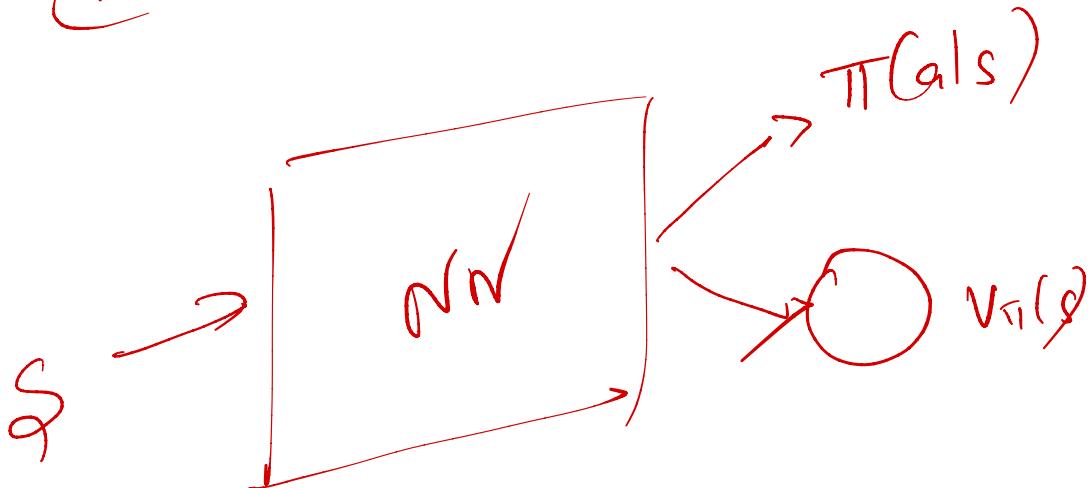
Importance weight

A3C

Asynchronous Advantage

Actor-Critic

(Mnih et al. 2016)



$$\theta_{t+1} = \theta_t + \alpha (G_{t:t+1} - \hat{v}(s; \omega)) \nabla \ln \pi$$

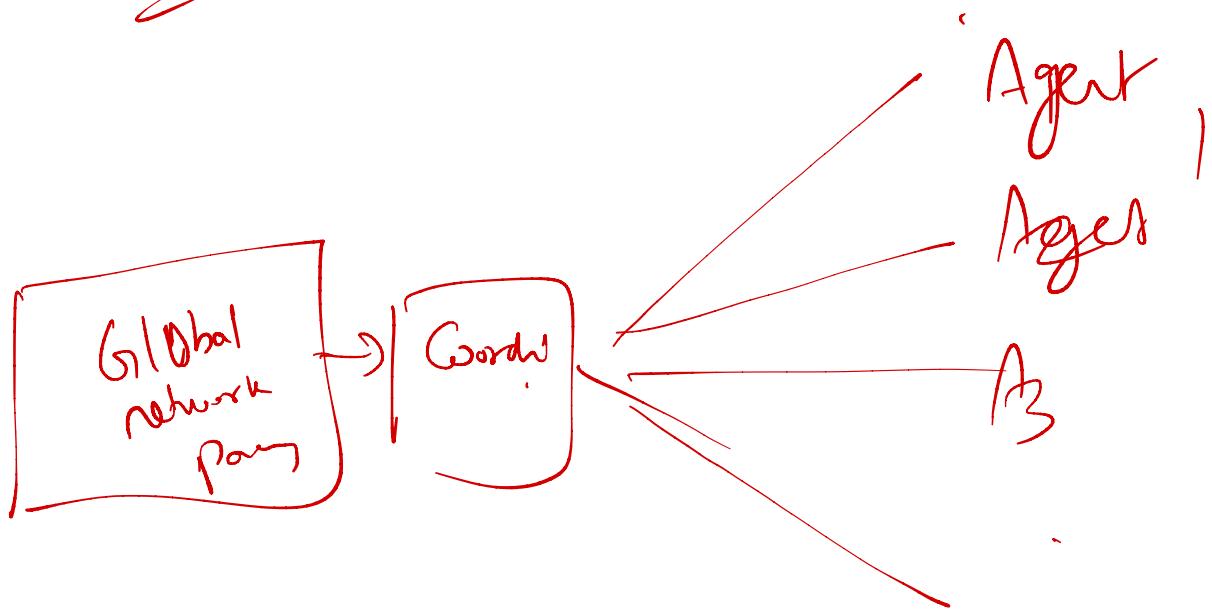
Advantage

$$(E(q(s,a)) - E(v_\pi(s)))$$

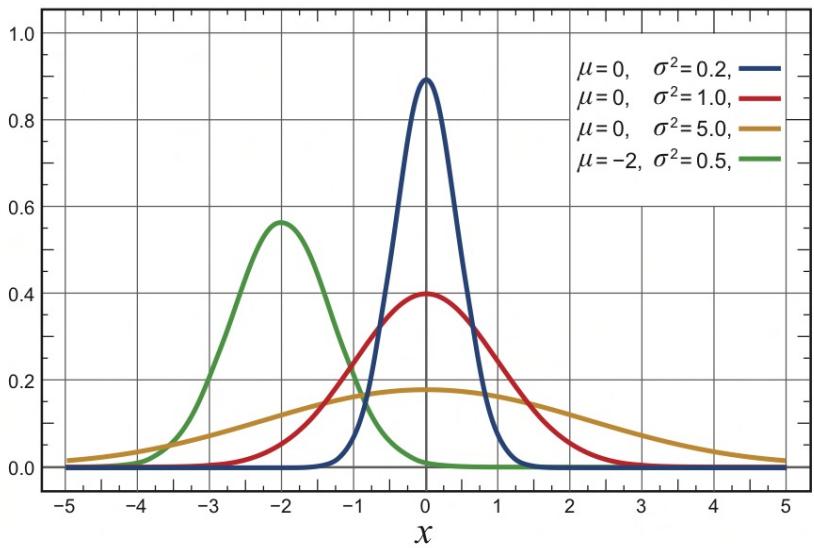
distributed training



A2C



Continuous actions:



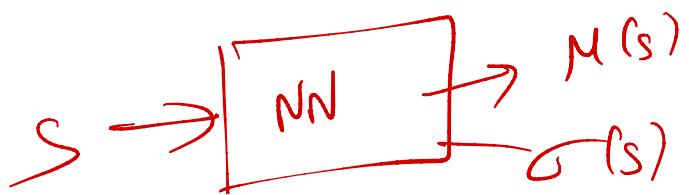
$$P(a) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(a-\mu)^2}{2\sigma^2}\right)$$

$$\pi(a|s, \theta) = \frac{1}{\sigma(s; \theta)\sqrt{2\pi}} \exp\left(-\frac{(a - \mu(s; \theta))^2}{2\sigma(s; \theta)^2}\right)$$

$$\theta = [\theta_\mu, \theta_\sigma]$$

$$\mu(s; \theta) = \theta_\mu^T x_\mu(s)$$

$$\sigma(s; \theta) = \exp(\theta_\sigma^T x_\sigma(s))$$



Policy Gradient methods

REINFORCE

" with baseline

Actor-Critic -

off-policy P.G.
Corollary: Actor / Arg reward

Continuous action space

Algorithm S3 Asynchronous advantage actor-critic - pseudocode for each actor-learner thread.

// Assume global shared parameter vectors θ and θ_v and global shared counter $T = 0$
// Assume thread-specific parameter vectors θ' and θ'_v
Initialize thread step counter $t \leftarrow 1$
repeat
 Reset gradients: $d\theta \leftarrow 0$ and $d\theta_v \leftarrow 0$.
 Synchronize thread-specific parameters $\theta' = \theta$ and $\theta'_v = \theta_v$
 $t_{start} = t$
 Get state s_t
 repeat
 Perform a_t according to policy $\pi(a_t|s_t; \theta')$
 Receive reward r_t and new state s_{t+1}
 $t \leftarrow t + 1$
 $T \leftarrow T + 1$
 until terminal s_t **or** $t - t_{start} == t_{max}$
 $R = \begin{cases} 0 & \text{for terminal } s_t \\ V(s_t, \theta'_v) & \text{for non-terminal } s_t // \text{Bootstrap from last state} \end{cases}$
 for $i \in \{t - 1, \dots, t_{start}\}$ **do**
 $R \leftarrow r_i + \gamma R$
 Accumulate gradients wrt θ' : $d\theta \leftarrow d\theta + \nabla_{\theta'} \log \pi(a_i|s_i; \theta')(R - V(s_i; \theta'_v))$
 Accumulate gradients wrt θ'_v : $d\theta_v \leftarrow d\theta_v + \partial (R - V(s_i; \theta'_v))^2 / \partial \theta'_v$
 end for
 Perform asynchronous update of θ using $d\theta$ and of θ_v using $d\theta_v$.
until $T > T_{max}$
