Computer Engineering Department

# Policy-based Theoretical Guarantees

**Mohammad Hossein Rohban, Ph.D.**

Spring 2025

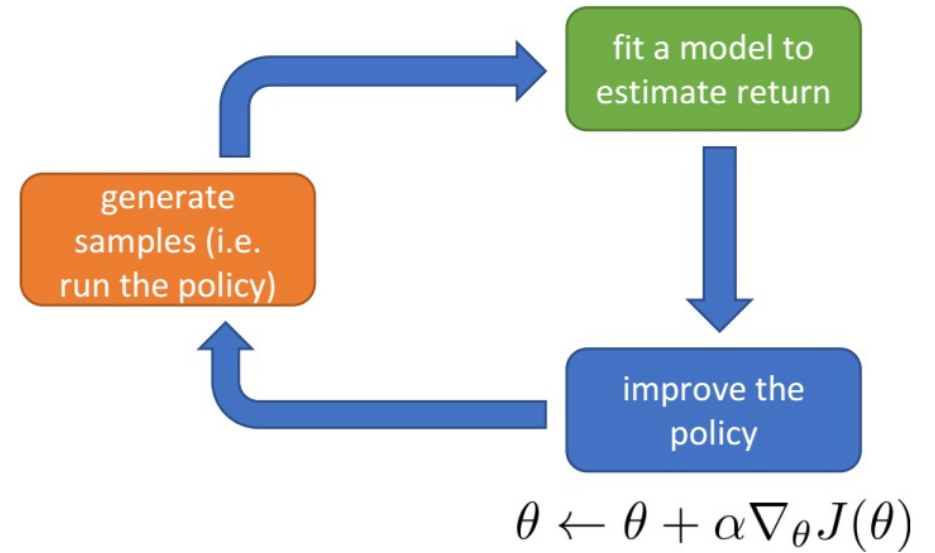Courtesy: Most of slides are adopted from the RL course at Berkeley.

REINFORCE algorithm:

1. sample $\{\tau^i\}$ from $\pi_\theta(\mathbf{a}_t|\mathbf{s}_t)$ (run the policy)
2. $\nabla_\theta J(\theta) \approx \sum_i \left( \sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_t^i|\mathbf{s}_t^i) \left( \sum_{t'=t}^T r(\mathbf{s}_{t'}^i, \mathbf{a}_{t'}^i) \right) \right)$
3. $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t}|\mathbf{s}_{i,t}) \hat{Q}_{i,t}^\pi$$

"reward to go"

$$\hat{Q}^\pi(\mathbf{x}_t, \mathbf{u}_t) = \sum_{t'=t}^T r(\mathbf{x}_{t'}, \mathbf{u}_{t'})$$

fit a model to estimate return

generate samples (i.e. run the policy)

improve the policy

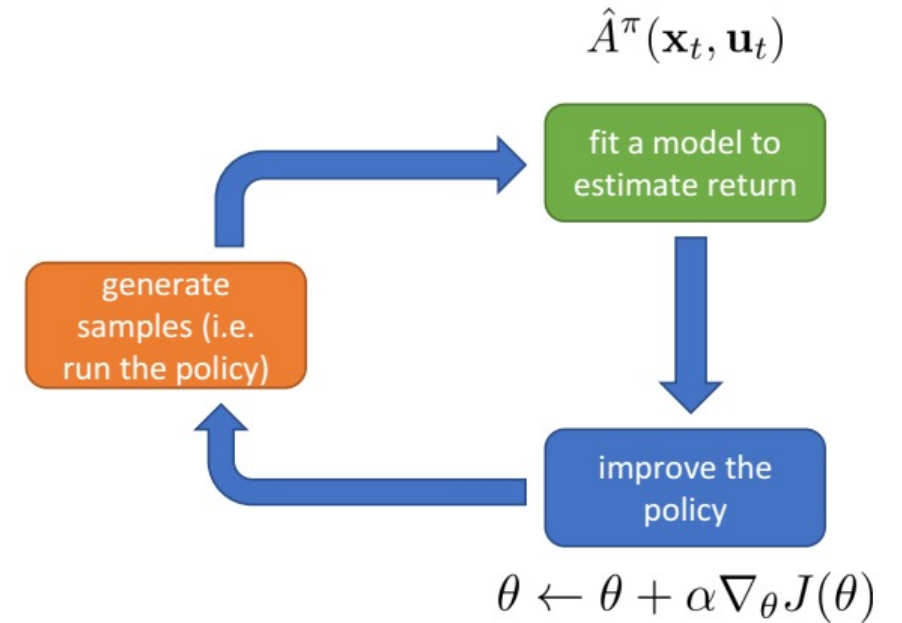$$\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$$

# Policy Gradient as Policy Iteration

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^{N} \sum_{t=1}^{T} \nabla_\theta \log \pi_\theta(\mathbf{a}_{i,t}|\mathbf{s}_{i,t}) \hat{A}_{i,t}^\pi$$

$$\hat{A}^\pi(\mathbf{x}_t, \mathbf{u}_t)$$

main steps of policy gradient algorithm:

1. Estimate $\hat{A}^\pi(\mathbf{s}_t, \mathbf{a}_t)$ for current policy $\pi$
2. Use $\hat{A}^\pi(\mathbf{s}_t, \mathbf{a}_t)$ to get *improved* policy $\pi'$

generate samples (i.e. run the policy)

fit a model to estimate return

improve the policy

$$\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$$

Familiar to policy iteration algorithm:

1. evaluate $A^\pi(\mathbf{s}, \mathbf{a})$ $\rightarrow$ $Q(s,a) - V(s)$
2. set $\pi \leftarrow \pi'$

$$J(\theta) = E_{\tau \sim p_\theta(\tau)} \left[ \sum_t \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) \right]$$

claim: $J(\theta') - J(\theta) = E_{\tau \sim p_{\theta'}(\tau)} \left[ \sum_t \gamma^t A^{\pi_\theta}(\mathbf{s}_t, \mathbf{a}_t) \right]$

could be interpreted as policy improvement!

claim: $J(\theta') - J(\theta) = E_{\tau \sim p_{\theta'}(\tau)} \left[ \sum_t \gamma^t A^{\pi_\theta}(\mathbf{s}_t, \mathbf{a}_t) \right]$

proof: 
$$J(\theta') - J(\theta) = J(\theta') - E_{\mathbf{s}_0 \sim p(\mathbf{s}_0)} \left[ V^{\pi_\theta}(\mathbf{s}_0) \right]$$

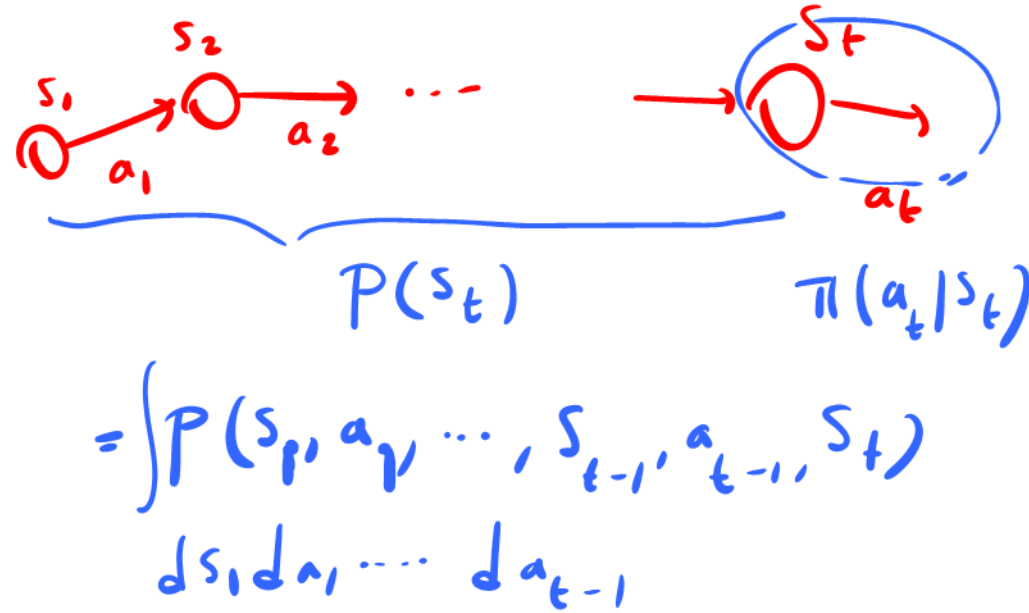$$= J(\theta') - E_{\tau \sim p_{\theta'}(\tau)} \left[ V^{\pi_\theta}(\mathbf{s}_0) \right]$$

$$= J(\theta') - E_{\tau \sim p_{\theta'}(\tau)} \left[ \sum_{t=0}^{\infty} \gamma^t V^{\pi_\theta}(\mathbf{s}_t) - \sum_{t=1}^{\infty} \gamma^t V^{\pi_\theta}(\mathbf{s}_t) \right]$$

$$= J(\theta') + E_{\tau \sim p_{\theta'}(\tau)} \left[ \sum_{t=0}^{\infty} \gamma^t (\gamma V^{\pi_\theta}(\mathbf{s}_{t+1}) - V^{\pi_\theta}(\mathbf{s}_t)) \right]$$

$$= E_{\tau \sim p_{\theta'}(\tau)} \left[ \sum_{t=0}^{\infty} \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) \right] + E_{\tau \sim p_{\theta'}(\tau)} \left[ \sum_{t=0}^{\infty} \gamma^t (\gamma V^{\pi_\theta}(\mathbf{s}_{t+1}) - V^{\pi_\theta}(\mathbf{s}_t)) \right]$$

$$= E_{\tau \sim p_{\theta'}(\tau)} \left[ \sum_{t=0}^{\infty} \gamma^t (r(\mathbf{s}_t, \mathbf{a}_t) + \gamma V^{\pi_\theta}(\mathbf{s}_{t+1}) - V^{\pi_\theta}(\mathbf{s}_t)) \right]$$

$$= E_{\tau \sim p_{\theta'}(\tau)} \left[ \sum_{t=0}^{\infty} \gamma^t A^{\pi_\theta}(\mathbf{s}_t, \mathbf{a}_t) \right]$$

$$P(s_t) = \int P(s_1, a_1, \cdots, s_{t-1}, a_{t-1}, s_t) \, ds_1 \, da_1 \cdots da_{t-1}$$

$$\pi(a_t | s_t)$$

$$J(\theta') - J(\theta) = E_{\tau \sim p_{\theta'}(\tau)} \left[ \sum_t \gamma^t A^{\pi_\theta}(\mathbf{s}_t, \mathbf{a}_t) \right]$$

expectation under $\pi_{\theta'}$          advantage under $\pi_\theta$

importance sampling

$$E_{x \sim p(x)}[f(x)] = \int p(x) f(x) dx$$
$$= \int \frac{q(x)}{q(x)} p(x) f(x) dx$$
$$= \int q(x) \frac{p(x)}{q(x)} f(x) dx$$
$$= E_{x \sim q(x)} \left[ \frac{p(x)}{q(x)} f(x) \right]$$

$$E_{\tau \sim p_{\theta'}(\tau)} \left[ \sum_t \gamma^t A^{\pi_\theta}(\mathbf{s}_t, \mathbf{a}_t) \right] = \sum_t E_{\mathbf{s}_t \sim p_{\theta'}(\mathbf{s}_t)} \left[ E_{\mathbf{a}_t \sim \pi_{\theta'}(\mathbf{a}_t|\mathbf{s}_t)} \left[ \gamma^t A^{\pi_\theta}(\mathbf{s}_t, \mathbf{a}_t) \right] \right]$$

$$= \sum_t E_{\mathbf{s}_t \sim p_{\theta'}(\mathbf{s}_t)} \left[ E_{\mathbf{a}_t \sim \pi_\theta(\mathbf{a}_t|\mathbf{s}_t)} \left[ \frac{\pi_{\theta'}(\mathbf{a}_t|\mathbf{s}_t)}{\pi_\theta(\mathbf{a}_t|\mathbf{s}_t)} \gamma^t A^{\pi_\theta}(\mathbf{s}_t, \mathbf{a}_t) \right] \right]$$

is it OK to use $p_\theta(\mathbf{s}_t)$ instead?

Can we ignore distribution mismatch?

$$\pi_\theta \approx \pi_{\theta'}$$

$$J(\theta') - J(\theta) =$$

$$\sum_t E_{\mathbf{s}_t \sim p_{\theta'}(\mathbf{s}_t)} \left[ E_{\mathbf{a}_t \sim \pi_\theta(\mathbf{a}_t | \mathbf{s}_t)} \left[ \frac{\pi_{\theta'}(\mathbf{a}_t | \mathbf{s}_t)}{\pi_\theta(\mathbf{a}_t | \mathbf{s}_t)} \gamma^t A^{\pi_\theta}(\mathbf{s}_t, \mathbf{a}_t) \right] \right] \overset{?}{\approx} \sum_t E_{\mathbf{s}_t \sim p_\theta(\mathbf{s}_t)} \left[ E_{\mathbf{a}_t \sim \pi_\theta(\mathbf{a}_t | \mathbf{s}_t)} \left[ \frac{\pi_{\theta'}(\mathbf{a}_t | \mathbf{s}_t)}{\pi_\theta(\mathbf{a}_t | \mathbf{s}_t)} \gamma^t A^{\pi_\theta}(\mathbf{s}_t, \mathbf{a}_t) \right] \right]$$

marg.

$$\pi(a_t | s_t)$$

**why do we want this to be true?**

$$P(s_t) \, P(a_t | s_t) = P(s_t, a_t)$$

$$\bar{A}(\theta')$$

$$J(\theta') - J(\theta) \approx \bar{A}(\theta') \quad \Rightarrow \quad \theta' \leftarrow \arg\max_{\theta'} \bar{A}(\theta)$$

2. Use $\hat{A}^\pi(\mathbf{s}_t, \mathbf{a}_t)$ to get *improved* policy $\pi'$

**is it true? and when?**

$p_\theta(\mathbf{s}_t)$ is *close* to $p_{\theta'}(\mathbf{s}_t)$ when $\pi_\theta$ is *close* to $\pi_{\theta'}$

Claim: $p_\theta(\mathbf{s}_t)$ is *close* to $p_{\theta'}(\mathbf{s}_t)$ when $\pi_\theta$ is *close* to $\pi_{\theta'}$

Simple case: assume $\pi_\theta$ is a *deterministic* policy $\mathbf{a}_t = \pi_\theta(\mathbf{s}_t)$

$\pi_{\theta'}$ is *close* to $\pi_\theta$ if $\pi_{\theta'}(\mathbf{a}_t \neq \pi_\theta(\mathbf{s}_t)|\mathbf{s}_t) \leq \epsilon$

$p_\theta(s_t) = (1-\epsilon)^t p_\theta(s_t) + (1-(1-\epsilon)^t) p_\theta(s_t)$

$$p_{\theta'}(\mathbf{s}_t) = (1-\epsilon)^t p_\theta(\mathbf{s}_t) + (1-(1-\epsilon)^t)) p_{\text{mistake}}(\mathbf{s}_t)$$

probability we made no mistakes          some *other* distribution

$$|p_{\theta'}(\mathbf{s}_t) - p_\theta(\mathbf{s}_t)| = (1-(1-\epsilon)^t)|p_{\text{mistake}}(\mathbf{s}_t) - p_\theta(\mathbf{s}_t)| \leq 2(1-(1-\epsilon)^t)$$

useful identity: $(1-\epsilon)^t \geq 1 - \epsilon t$ for $\epsilon \in [0,1]$          $\leq 2\epsilon t$

**not a great bound, but a bound!**

$\pi_{\theta'}$          95%.          $\pi_\theta(s) = 3$

$S$          1%.   1%   2%.   1%.

seem familiar?

$P(A) = \sum_i P(A \cap B_i)$

① $P(A|B_i) P(B_i)$

$S$

②

# Bounding the distribution change

$$\sum_y P(x,y) = P_X(x)$$

Claim: $p_\theta(\mathbf{s}_t)$ is *close* to $p_{\theta'}(\mathbf{s}_t)$ when $\pi_\theta$ is *close* to $\pi_{\theta'}$

$$P_X(x)$$
$$P_Y(y)$$

$$P(x,y) \text{ s.t.}$$

General case: assume $\pi_\theta$ is an arbitrary distribution

$$\sum_x P(x,y) = P_Y(y)$$

$\pi_{\theta'}$ is *close* to $\pi_\theta$ if $|\pi_{\theta'}(\mathbf{a}_t|\mathbf{s}_t) - \pi_\theta(\mathbf{a}_t|\mathbf{s}_t)| \le \epsilon$ for all $\mathbf{s}_t$

$$P(X = Y) = 1 - \epsilon$$

Useful lemma: if $|p_X(x) - p_Y(x)| = \epsilon$, exists $p(x, y)$ such that $p(x) = p_X(x)$ and $p(y) = p_Y(y)$ and $p(x = y) = 1 - \epsilon$

$\Rightarrow p_X(x)$ "agrees" with $p_Y(y)$ with probability $\epsilon$

$\Rightarrow \pi_{\theta'}(\mathbf{a}_t|\mathbf{s}_t)$ takes a different action than $\pi_\theta(\mathbf{a}_t|\mathbf{s}_t)$ with probability at most $\epsilon$
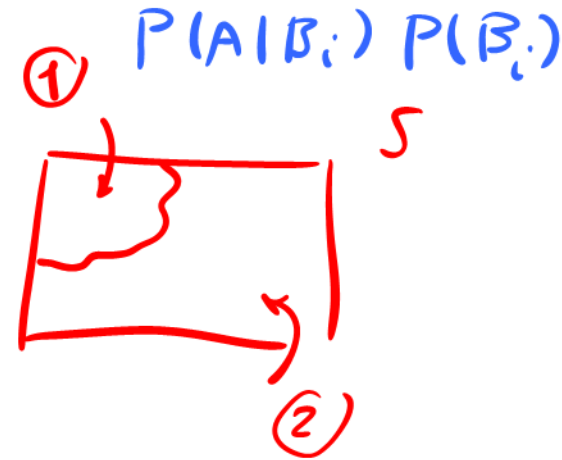
$$|p_{\theta'}(\mathbf{s}_t) - p_\theta(\mathbf{s}_t)| = (1 - (1 - \epsilon)^t)|p_{\text{mistake}}(\mathbf{s}_t) - p_\theta(\mathbf{s}_t)| \le 2(1 - (1 - \epsilon)^t)$$
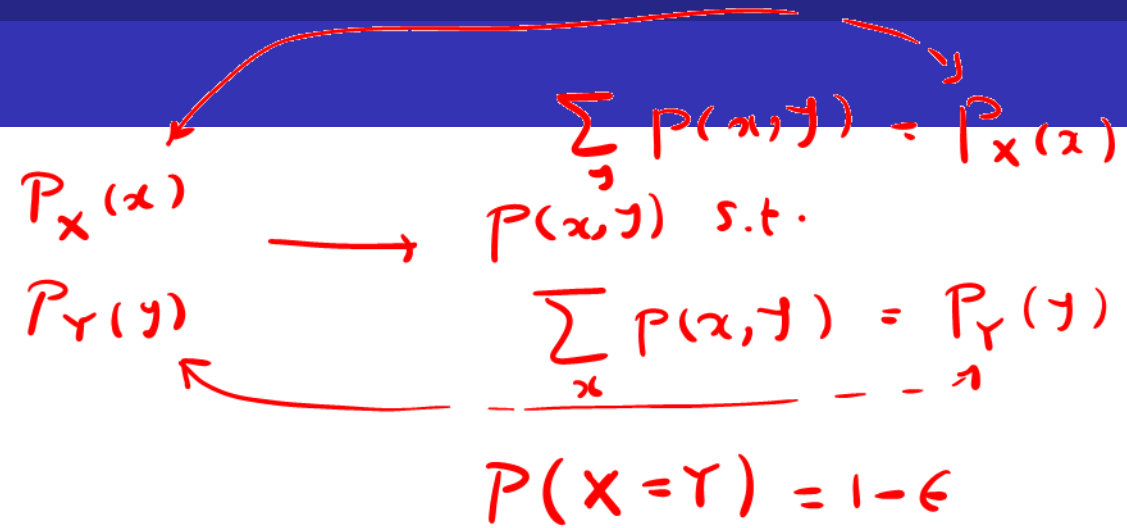
$$\le 2\epsilon t$$

$\pi_{\theta'}$ is *close* to $\pi_\theta$ if $|\pi_{\theta'}(\mathbf{a}_t|\mathbf{s}_t) - \pi_\theta(\mathbf{a}_t|\mathbf{s}_t)| \le \epsilon$ for all $\mathbf{s}_t$

$|p_{\theta'}(\mathbf{s}_t) - p_\theta(\mathbf{s}_t)| \le 2\epsilon t$

$$\sum_s \left[ P_{\theta'}(s) - P_\theta(s) \right] f(s)$$

$$\rightarrow P_{\theta'} + P_\theta - P_\theta$$

$$E_{p_{\theta'}(\mathbf{s}_t)}[f(\mathbf{s}_t)] = \sum_{\mathbf{s}_t} p_{\theta'}(\mathbf{s}_t)f(\mathbf{s}_t) \ge \sum_{\mathbf{s}_t} p_\theta(\mathbf{s}_t)f(\mathbf{s}_t) - |p_\theta(\mathbf{s}_t) - p_{\theta'}(\mathbf{s}_t)| \max_{\mathbf{s}_t} f(\mathbf{s}_t)$$

$$\ge E_{p_\theta(\mathbf{s}_t)}[f(\mathbf{s}_t)] - 2\epsilon t \max_{\mathbf{s}_t} f(\mathbf{s}_t)$$

$$\sum_t E_{\mathbf{s}_t \sim p_{\theta'}(\mathbf{s}_t)} \left[ E_{\mathbf{a}_t \sim \pi_\theta(\mathbf{a}_t|\mathbf{s}_t)} \left[ \frac{\pi_{\theta'}(\mathbf{a}_t|\mathbf{s}_t)}{\pi_\theta(\mathbf{a}_t|\mathbf{s}_t)} \gamma^t A^{\pi_\theta}(\mathbf{s}_t, \mathbf{a}_t) \right] \right] \ge$$

$$O(T r_{\max}) \text{ or } O\left(\frac{r_{\max}}{1-\gamma}\right)$$

$$\sum_t E_{\mathbf{s}_t \sim p_\theta(\mathbf{s}_t)} \left[ E_{\mathbf{a}_t \sim \pi_\theta(\mathbf{a}_t|\mathbf{s}_t)} \left[ \frac{\pi_{\theta'}(\mathbf{a}_t|\mathbf{s}_t)}{\pi_\theta(\mathbf{a}_t|\mathbf{s}_t)} \gamma^t A^{\pi_\theta}(\mathbf{s}_t, \mathbf{a}_t) \right] \right] - \sum_t 2\epsilon t C$$

maximizing this maximizes a bound on the thing we want!

1. **Q-function update**
   Update Q-function to evaluate current policy:

   $$Q(\mathbf{s}, \mathbf{a}) \leftarrow r(\mathbf{s}, \mathbf{a}) + \mathbb{E}_{\mathbf{s}' \sim p_\mathbf{s}, \ \mathbf{a}' \sim \pi} \left[ Q(\mathbf{s}', \mathbf{a}') - \log \pi(\mathbf{a}'|\mathbf{s}') \right]$$

   This converges to $Q^\pi$.

2. **Update policy**
   Update the policy with gradient of information projection:

   $$\pi_{\text{new}} = \arg\min_{\pi'} D_{\text{KL}} \left( \pi'_\theta(\cdot|\mathbf{s}) \ \middle\| \ \frac{1}{Z} \exp Q^{\pi_{\text{old}}}(\mathbf{s}, \cdot) \right)$$

   *(annotations: θ, NN)*

   In practice, only take one gradient step on this objective

3. **Interact with the world, collect more data**

Haarnoja, et al. **Soft Actor-Critic Algorithms and Applications**. '18

# Soft actor-critic

**Algorithm 1** Soft Actor-Critic

**Inputs**: The learning rates, $\lambda_\pi$, $\lambda_Q$, and $\lambda_V$ for functions $\pi_\theta$, $Q_w$, and $V_\psi$ respectively; the weighting factor $\tau$ for exponential moving average.

1: Initialize parameters $\theta$, $w$, $\psi$, and $\bar{\psi}$.
2: **for** each iteration **do**
3:    *(In practice, a combination of a single environment step and multiple gradient steps is found to work best.)*
4:        **for** each environment setup **do**
5:            $a_t \sim \pi_\theta(a_t|s_t)$
6:            $s_{t+1} \sim \rho_\pi(s_{t+1}|s_t, a_t)$
7:            $\mathcal{D} \leftarrow \mathcal{D} \cup \{(s_t, a_t, r(s_t, a_t), s_{t+1}\}$
8:        **for** each gradient update step **do**
9:            $\psi \leftarrow \psi - \lambda_V \nabla_\psi J_V(\psi).$
10:            $w \leftarrow w - \lambda_Q \nabla_w J_Q(w).$
11:            $\theta \leftarrow \theta - \lambda_\pi \nabla_\theta J_\pi(\theta).$
12:            $\bar{\psi} \leftarrow \tau\psi + (1-\tau)\bar{\psi}.$

# Loss functions

$$J_V(\psi) = \mathbb{E}_{\mathbf{s}_t \sim \mathcal{D}} \left[ \tfrac{1}{2} \left( V_\psi(\mathbf{s}_t) - \mathbb{E}_{\mathbf{a}_t \sim \pi_\phi} \left[ Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \log \pi_\phi(\mathbf{a}_t|\mathbf{s}_t) \right] \right)^2 \right]$$

(5)

$$J_Q(\theta) = \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \mathcal{D}} \left[ \tfrac{1}{2} \left( Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) \right)^2 \right],$$

$$\sum_{s_t \in D} \left[ \sum_{a_t} \pi_\phi(a_t \mid s_t) \, \log Z_\theta(s_t) \quad \cdots \quad \right]$$

(7)

with

$$p(s_{t+1} \mid s_t, a_t)$$

$$\hat{Q}(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \, \mathbb{E}_{\mathbf{s}_{t+1} \sim p} \left[ V_{\bar{\psi}}(\mathbf{s}_{t+1}) \right],$$

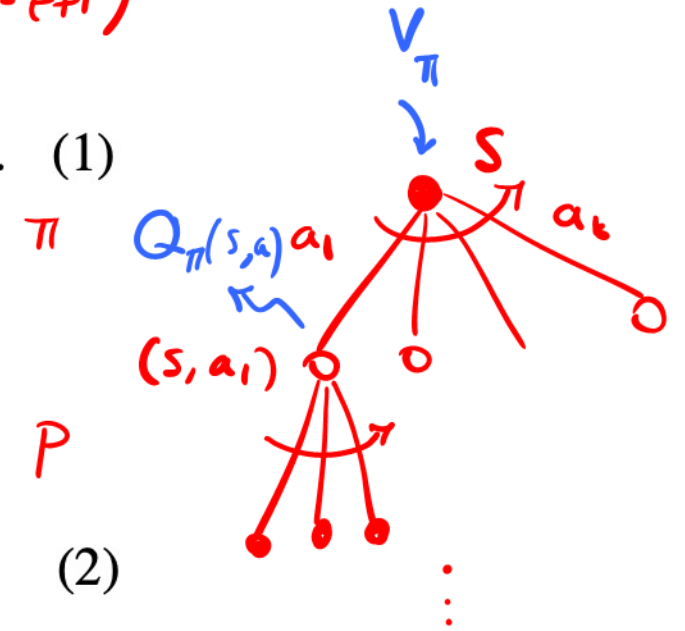(8)

$$J_\pi(\phi) = \mathbb{E}_{\mathbf{s}_t \sim \mathcal{D}} \left[ \mathrm{D}_{\mathrm{KL}} \left( \pi_\phi(\cdot|\mathbf{s}_t) \, \Big\| \, \frac{\exp\left(Q_\theta(\mathbf{s}_t, \cdot)\right)}{Z_\theta(\mathbf{s}_t)} \right) \right].$$

(10)

$$\sum_{s_{t+1},\, a_{t+1}} \log \pi(a_{t+1}|s_{t+1}) \cdot P(s_{t+1}|s_t, a_t)\, \pi(a_{t+1}|s_{t+1})$$

$$J(\pi) = \sum_{t=0}^{T} \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \rho_\pi} \left[ r(\mathbf{s}_t, \mathbf{a}_t) + \alpha \mathcal{H}(\pi(\cdot|\mathbf{s}_t)) \right]. \quad (1)$$

$V_\pi$

$\pi$ $Q_\pi(s,a)\, a_1$

$(s, a_1)$

$r'(s_t, a_t)$

$P$

$$\mathcal{T}^\pi Q(\mathbf{s}_t, \mathbf{a}_t) \triangleq r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \mathbb{E}_{\mathbf{s}_{t+1} \sim p}\left[ V(\mathbf{s}_{t+1}) \right], \quad (2)$$

where

$$P(s_{t+1}|s_t, a_t) \cdot \pi(a_{t+1}|s_{t+1})$$

$$V(\mathbf{s}_t) = \mathbb{E}_{\mathbf{a}_t \sim \pi}\left[ Q(\mathbf{s}_t, \mathbf{a}_t) - \log \pi(\mathbf{a}_t|\mathbf{s}_t) \right] \quad (3)$$

$$\mathbb{E}_{s_{t+1}, \pi}\left[ Q(s_t, a_t) - \log \pi(a_t|s_{t+1}) \right]$$

$$g(s_t, a_t)$$

$$\mathbb{E}_{s_{t+1}, \pi}\left[ Q(s_t, a_t) - const(s_{t+1}, a_{t+1}) \right]$$

# Soft Policy Evaluation

**Lemma 1** (Soft Policy Evaluation). *Consider the soft Bellman backup operator $\mathcal{T}^\pi$ in* Equation 2 *and a mapping $Q^0 : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ with $|\mathcal{A}| < \infty$, and define $Q^{k+1} = \mathcal{T}^\pi Q^k$. Then the sequence $Q^k$ will converge to the soft Q-value of $\pi$ as $k \to \infty$.*

# Soft Policy Evaluation

**Lemma 1** (Soft Policy Evaluation). *Consider the soft Bellman backup operator $\mathcal{T}^\pi$ in Equation 2 and a mapping $Q^0 : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ with $|\mathcal{A}| < \infty$, and define $Q^{k+1} = \mathcal{T}^\pi Q^k$. Then the sequence $Q^k$ will converge to the soft Q-value of $\pi$ as $k \to \infty$.*

*Proof.* Define the entropy augmented reward as $r_\pi(\mathbf{s}_t, \mathbf{a}_t) \triangleq r(\mathbf{s}_t, \mathbf{a}_t) + \mathbb{E}_{\mathbf{s}_{t+1} \sim p} [\mathcal{H}(\pi(\cdot | \mathbf{s}_{t+1}))]$ and rewrite the update rule as

$$Q(\mathbf{s}_t, \mathbf{a}_t) \leftarrow r_\pi(\mathbf{s}_t, \mathbf{a}_t) + \gamma \mathbb{E}_{\mathbf{s}_{t+1} \sim p, \mathbf{a}_{t+1} \sim \pi} [Q(\mathbf{s}_{t+1}, \mathbf{a}_{t+1})] \tag{15}$$

and apply the standard convergence results for policy evaluation (Sutton & Barto, 1998). The assumption $|\mathcal{A}| < \infty$ is required to guarantee that the entropy augmented reward is bounded. $\square$

# Soft Policy Improvement

$$\pi_{\text{new}} = \arg\min_{\pi' \in \Pi} \mathrm{D}_{\mathrm{KL}}\left(\pi'(\cdot|\mathbf{s}_t) \,\middle\|\, \frac{\exp\left(Q^{\pi_{\text{old}}}(\mathbf{s}_t, \cdot)\right)}{Z^{\pi_{\text{old}}}(\mathbf{s}_t)}\right).$$

$$(4)$$

**Lemma 2** (Soft Policy Improvement). *Let $\pi_{\text{old}} \in \Pi$ and let $\pi_{\text{new}}$ be the optimizer of the minimization problem defined in Equation 4. Then $Q^{\pi_{\text{new}}}(\mathbf{s}_t, \mathbf{a}_t) \geq Q^{\pi_{\text{old}}}(\mathbf{s}_t, \mathbf{a}_t)$ for all $(\mathbf{s}_t, \mathbf{a}_t) \in \mathcal{S} \times \mathcal{A}$ with $|\mathcal{A}| < \infty$.*

$$D_{KL}(P \| q) = \sum_x P(x) \log \frac{P(x)}{q(x)} = \sum P(x) \log P(x) - P(x) \log q(x)$$

**Lemma 2** (Soft Policy Improvement). *Let $\pi_{old} \in \Pi$ and let $\pi_{new}$ be the optimizer of the minimization problem defined in Equation 4. Then $Q^{\pi_{new}}(\mathbf{s}_t, \mathbf{a}_t) \geq Q^{\pi_{old}}(\mathbf{s}_t, \mathbf{a}_t)$ for all $(\mathbf{s}_t, \mathbf{a}_t) \in \mathcal{S} \times \mathcal{A}$ with $|\mathcal{A}| < \infty$.*

*Proof.* Let $\pi_{old} \in \Pi$ and let $Q^{\pi_{old}}$ and $V^{\pi_{old}}$ be the corresponding soft state-action value and soft state value, and let $\pi_{new}$ be defined as

$$\sum_{a_t} \pi'(a_t|s_t) \log \pi' - \pi' \log \exp(\cdots)$$

$$\pi_{new}(\cdot|\mathbf{s}_t) = \arg\min_{\pi' \in \Pi} D_{KL}\left(\pi'(\cdot|\mathbf{s}_t) \| \exp\left(Q^{\pi_{old}}(\mathbf{s}_t, \cdot) - \log Z^{\pi_{old}}(\mathbf{s}_t)\right)\right)$$

$$\mathbb{E}_{a \sim \pi'}\left[\log \pi' - Q^{\pi_{old}} + \log Z\right] = \arg\min_{\pi' \in \Pi} J_{\pi_{old}}(\pi'(\cdot|\mathbf{s}_t)). \qquad \sum \pi' \log \pi' - \pi'\left[Q^{\pi_{old}} - \log Z\right] \quad (16)$$

It must be the case that $J_{\pi_{old}}(\pi_{new}(\cdot|\mathbf{s}_t)) \leq J_{\pi_{old}}(\pi_{old}(\cdot|\mathbf{s}_t))$, since we can always choose $\pi_{new} = \pi_{old} \in \Pi$. Hence

$$\mathbb{E}_{\mathbf{a}_t \sim \pi_{new}}\left[\log \pi_{new}(\mathbf{a}_t|\mathbf{s}_t) - Q^{\pi_{old}}(\mathbf{s}_t, \mathbf{a}_t) + \log Z^{\pi_{old}}(\mathbf{s}_t)\right] \leq \mathbb{E}_{\mathbf{a}_t \sim \pi_{old}}\left[\log \pi_{old}(\mathbf{a}_t|\mathbf{s}_t) - Q^{\pi_{old}}(\mathbf{s}_t, \mathbf{a}_t) + \log Z^{\pi_{old}}(\mathbf{s}_t)\right],$$
$$(17)$$

$$-V^{\pi_{old}}(s_t)$$

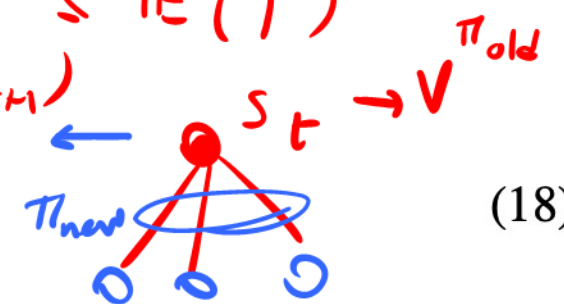$$\mathbb{E}_{a_{t+1} \sim \pi_{new}}\left[ Q^{\pi_{old}}(s_{t+1}, a_{t+1}) - \log \pi_{new}(a_{t+1}|s_{t+1}) \right] \geq V^{\pi_{old}}(s_{t+1})$$

$$X \leq Y \longrightarrow \mathbb{E}(X) \leq \mathbb{E}(Y)$$

and since partition function $Z^{\pi_{old}}$ depends only on the state, the inequality reduces to

$$\mathbb{E}_{\mathbf{a}_t \sim \pi_{new}}\left[ Q^{\pi_{old}}(\mathbf{s}_t, \mathbf{a}_t) - \log \pi_{new}(\mathbf{a}_t|\mathbf{s}_t) \right] \geq V^{\pi_{old}}(\mathbf{s}_t). \tag{18}$$

Next, consider the soft Bellman equation:

$$Q^{\pi_{old}}(\mathbf{s}_t, \mathbf{a}_t) = r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \, \mathbb{E}_{\mathbf{s}_{t+1} \sim p}\left[ V^{\pi_{old}}(\mathbf{s}_{t+1}) \right]$$

$$\leq r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \, \mathbb{E}_{\mathbf{s}_{t+1} \sim p}\left[ \mathbb{E}_{\mathbf{a}_{t+1} \sim \pi_{new}}\left[ Q^{\pi_{old}}(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}) - \log \pi_{new}(\mathbf{a}_{t+1}|\mathbf{s}_{t+1}) \right] \right]$$

$$\vdots$$

$$\leq Q^{\pi_{new}}(\mathbf{s}_t, \mathbf{a}_t), \tag{19}$$

where we have repeatedly expanded $Q^{\pi_{old}}$ on the RHS by applying the soft Bellman equation and the bound in Equation 18. Convergence to $Q^{\pi_{new}}$ follows from Lemma 1. $\square$

**Theorem 1** (Soft Policy Iteration). *Repeated application of soft policy evaluation and soft policy improvement to any $\pi \in \Pi$ converges to a policy $\pi^*$ such that $Q^{\pi^*}(\mathbf{s}_t, \mathbf{a}_t) \geq Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t)$ for all $\pi \in \Pi$ and $(\mathbf{s}_t, \mathbf{a}_t) \in \mathcal{S} \times \mathcal{A}$, assuming $|\mathcal{A}| < \infty$.*

*Proof.* Let $\pi_i$ be the policy at iteration $i$. By Lemma 2, the sequence $Q^{\pi_i}$ is monotonically increasing. Since $Q^{\pi}$ is bounded above for $\pi \in \Pi$ (both the reward and entropy are bounded), the sequence converges to some $\pi^*$. We will still need to show that $\pi^*$ is indeed optimal. At convergence, it must be case that $J_{\pi^*}(\pi^*(\cdot|\mathbf{s}_t)) < J_{\pi^*}(\pi(\cdot|\mathbf{s}_t))$ for all $\pi \in \Pi$, $\pi \neq \pi^*$. Using the same iterative argument as in the proof of Lemma 2, we get $Q^{\pi^*}(\mathbf{s}_t, \mathbf{a}_t) > Q^{\pi}(\mathbf{s}_t, \mathbf{a}_t)$ for all $(\mathbf{s}_t, \mathbf{a}_t) \in \mathcal{S} \times \mathcal{A}$, that is, the soft value of any other policy in $\Pi$ is lower than that of the converged policy. Hence $\pi^*$ is optimal in $\Pi$. $\square$

Policy @
Convergence $\rightarrow$ $\pi^* = \underset{\pi' \in \Pi}{\arg\min} \; J_{\pi^*}(\pi') \quad \rightarrow \quad J_{\pi^*}(\pi^*) \leq J_{\pi^*}(\pi)$

$\pi_{new}$ $\quad$ $\pi' \in \Pi$ $\quad$ $\pi_{old}$ $\qquad\qquad$ $\pi^*$ $\qquad$ $\pi^*$

$$\underbrace{\mathbb{E}_{\pi^*}\left[ Q^{\pi^*}(s_t, a_t) - \lg \pi_*(a_t|s_t) \right]}_{V^{\pi^*}(s_t)} \geq \mathbb{E}_{\pi}\left[ Q^{\pi^*}(s_t, a_t) - \lg \pi \right]$$