Hill Climbing

$$E[R^{\pi_\theta}]$$

$\theta$

# Policy Search

## Hill Climbing

$$E[R^{\pi_\theta}]$$

$\theta$

## Genetic Search

$\theta_1$ | $\theta_1^1$ | $\theta_1^2$ | $\theta_1^3$ | $\theta_1^4$ |

$\theta_2$ | $\theta_2^1$ | $\theta_2^2$ | $\theta_2^3$ | $\theta_2^4$ |

$\vdots$

$\theta_N$ | $\theta_N^1$ | $\theta_N^2$ | $\theta_N^3$ | $\theta_N^4$ |

# Policy Search

## Hill Climbing

$$E[R^{\pi_\theta}]$$

$\theta$

## Genetic Search

$\theta_1$ | $\theta_1^1$ | $\theta_1^2$ | $\theta_1^3$ | $\theta_1^4$

$\theta_2$ | $\theta_2^1$ | $\theta_2^2$ | $\theta_2^3$ | $\theta_2^4$

$\vdots$

$\theta_N$ | $\theta_N^1$ | $\theta_N^2$ | $\theta_N^3$ | $\theta_N^4$

crossover

mutation

evaluation
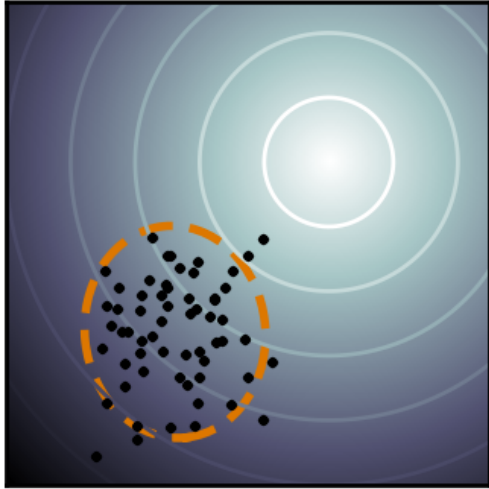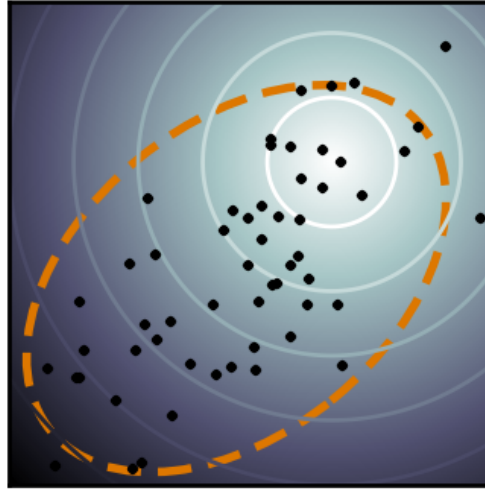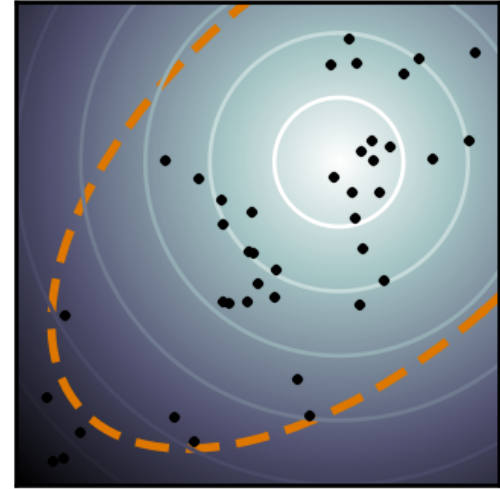
# Policy Search

## CMA-ES



Generation 1     Generation 2     Generation 3

Generation 4     Generation 5     Generation 6

# Gradient Bandits:

$$H_{t+1}(a) \doteq H_t(a) + \alpha \frac{\partial \mathbb{E}[R_t]}{\partial H_t(a)}$$

$$\frac{\partial \mathbb{E}[R_t]}{\partial H_t(a)} = \frac{\partial}{\partial H_t(a)} \left[ \sum_x \pi_t(x) q_*(x) \right]$$
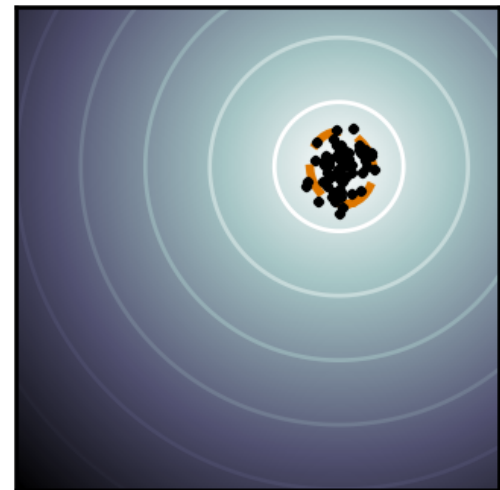
# Gradient Bandits:

$$H_{t+1}(a) \doteq H_t(a) + \alpha \frac{\partial \mathbb{E}[R_t]}{\partial H_t(a)}$$

$$\frac{\partial \mathbb{E}[R_t]}{\partial H_t(a)} = \frac{\partial}{\partial H_t(a)} \left[ \sum_x \pi_t(x) q_*(x) \right]$$

Just a scalar per arm. No states!

But in full RL case, policy influenes future states!

$$\nabla v_\pi(s) = \nabla \left[ \sum_a \pi(a|s) q_\pi(s,a) \right], \quad \text{for all } s \in \mathcal{S} \qquad \text{(Exercise 3.18)}$$

$$= \sum_a \left[ \nabla \pi(a|s) q_\pi(s,a) + \pi(a|s) \nabla q_\pi(s,a) \right] \quad \text{(product rule of calculus)}$$

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

$$\text{(Exercise 3.19 and Equation 3.2)}$$

$$= \sum_a \left[ \nabla \pi(a|s) q_\pi(s,a) + \pi(a|s) \sum_{s'} p(s'|s,a) \nabla v_\pi(s') \right] \qquad \text{(Eq. 3.4)}$$

$$= \sum_a \left[ \nabla \pi(a|s) q_\pi(s,a) + \pi(a|s) \sum_{s'} p(s'|s,a) \right. \qquad \text{(unrolling)}$$

$$\left. \sum_{a'} \left[ \nabla \pi(a'|s') q_\pi(s',a') + \pi(a'|s') \sum_{s''} p(s''|s',a') \nabla v_\pi(s'') \right] \right]$$

$$= \sum_{x \in \mathcal{S}} \sum_{k=0}^{\infty} \Pr(s \to x, k, \pi) \sum_a \nabla \pi(a|x) q_\pi(x,a),$$

after repeated unrolling, where $\Pr(s \to x, k, \pi)$ is the probability of transitioning from state $s$ to state $x$ in $k$ steps under policy $\pi$. It is then immediate that

$$\nabla J(\boldsymbol{\theta}) = \nabla v_\pi(s_0)$$

$$= \sum_s \left( \sum_{k=0}^{\infty} \Pr(s_0 \to s, k, \pi) \right) \sum_a \nabla \pi(a|s) q_\pi(s,a)$$

$$= \sum_s \eta(s) \sum_a \nabla \pi(a|s) q_\pi(s,a) \qquad \text{(box page 199)}$$

$$= \sum_{s'} \eta(s') \sum_s \frac{\eta(s)}{\sum_{s'} \eta(s')} \sum_a \nabla \pi(a|s) q_\pi(s,a)$$

$$= \sum_{s'} \eta(s') \sum_s \mu(s) \sum_a \nabla \pi(a|s) q_\pi(s,a) \qquad \text{(Eq. 9.3)}$$

$$\propto \sum_s \mu(s) \sum_a \nabla \pi(a|s) q_\pi(s,a) \qquad \text{(Q.E.D.)}$$

$$\nabla v_\pi(s) = \nabla \left[ \sum_a \pi(a|s) q_\pi(s,a) \right], \quad \text{for all } s \in \mathcal{S} \qquad \text{(Exercise 3.18)}$$

$$= \sum_a \left[ \nabla \pi(a|s) q_\pi(s,a) + \pi(a|s) \nabla q_\pi(s,a) \right] \quad \text{(product rule of calculus)}$$

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

$$\text{(Exercise 3.19 and Equation 3.2)}$$

$$= \sum_a \left[ \nabla \pi(a|s) q_\pi(s,a) + \pi(a|s) \sum_{s'} p(s'|s,a) \nabla v_\pi(s') \right] \qquad \text{(Eq. 3.4)} \longrightarrow$$

$$= \sum_a \left[ \nabla \pi(a|s) q_\pi(s,a) + \pi(a|s) \sum_{s'} p(s'|s,a) \right] \qquad \text{(unrolling)}$$

$$\sum_{a'} \left[ \nabla \pi(a'|s') q_\pi(s',a') + \pi(a'|s') \sum_{s''} p(s''|s',a') \nabla v_\pi(s'') \right] \Big]$$

$$= \sum_{x \in \mathcal{S}} \sum_{k=0}^\infty \Pr(s \to x, k, \pi) \sum_a \nabla \pi(a|x) q_\pi(x,a),$$

after repeated unrolling, where $\Pr(s \to x, k, \pi)$ is the probability of transitioning from state $s$ to state $x$ in $k$ steps under policy $\pi$. It is then immediate that

$$\nabla J(\boldsymbol{\theta}) = \nabla v_\pi(s_0)$$

$$= \sum_s \left( \sum_{k=0}^\infty \Pr(s_0 \to s, k, \pi) \right) \sum_a \nabla \pi(a|s) q_\pi(s,a)$$

$$= \sum_s \eta(s) \sum_a \nabla \pi(a|s) q_\pi(s,a) \qquad \text{(box page 199)}$$

$$= \sum_{s'} \eta(s') \sum_s \frac{\eta(s)}{\sum_{s'} \eta(s')} \sum_a \nabla \pi(a|s) q_\pi(s,a)$$

$$= \sum_{s'} \eta(s') \sum_s \mu(s) \sum_a \nabla \pi(a|s) q_\pi(s,a) \qquad \text{(Eq. 9.3)}$$

$$\propto \sum_s \mu(s) \sum_a \nabla \pi(a|s) q_\pi(s,a) \qquad \text{(Q.E.D.)}$$

*(handwritten annotations)*

Proof of policy gradient theorem

Marginalize R, push in gradient

Dynamics + Reward constant w.r.t. $\theta$

$$\nabla v_\pi(s) = \nabla\left[\sum_a \pi(a|s)q_\pi(s,a)\right], \quad \text{for all } s \in \mathcal{S} \qquad \text{(Exercise 3.18)}$$

$$= \sum_a \left[\nabla\pi(a|s)q_\pi(s,a) + \pi(a|s)\nabla q_\pi(s,a)\right] \quad \text{(product rule of calculus)}$$

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

$$\text{(Exercise 3.19 and Equation 3.2)}$$

$$= \sum_a \left[\nabla\pi(a|s)q_\pi(s,a) + \pi(a|s)\sum_{s'} p(s'|s,a)\nabla v_\pi(s')\right] \qquad \text{(Eq. 3.4)} \longrightarrow$$

$$= \sum_a \left[\nabla\pi(a|s)q_\pi(s,a) + \pi(a|s)\sum_{s'} p(s'|s,a) \right. \qquad \text{(unrolling)}$$

$$\left. \sum_{a'}\left[\nabla\pi(a'|s')q_\pi(s',a') + \pi(a'|s')\sum_{s''} p(s''|s',a')\nabla v_\pi(s'')\right]\right]$$

$$= \sum_{x\in\mathcal{S}}\sum_{k=0}^{\infty} \Pr(s\to x, k, \pi)\sum_a \nabla\pi(a|x)q_\pi(x,a), \qquad \cdots\longrightarrow$$

after repeated unrolling, where $\Pr(s\to x, k, \pi)$ is the probability of transitioning from state $s$ to state $x$ in $k$ steps under policy $\pi$. It is then immediate that

$$\nabla J(\boldsymbol{\theta}) = \nabla v_\pi(s_0)$$

$$= \sum_s \left(\sum_{k=0}^{\infty} \Pr(s_0\to s, k, \pi)\right)\sum_a \nabla\pi(a|s)q_\pi(s,a)$$

$$= \sum_s \eta(s)\sum_a \nabla\pi(a|s)q_\pi(s,a) \qquad \text{(box page 199)}$$

$$= \sum_{s'} \eta(s')\sum_s \frac{\eta(s)}{\sum_{s'}\eta(s')}\sum_a \nabla\pi(a|s)q_\pi(s,a)$$

$$= \sum_{s'} \eta(s')\sum_s \mu(s)\sum_a \nabla\pi(a|s)q_\pi(s,a) \qquad \text{(Eq. 9.3)}$$

$$\propto \sum_s \mu(s)\sum_a \nabla\pi(a|s)q_\pi(s,a) \qquad \text{(Q.E.D.)}$$

*Handwritten annotations:*

Proof of policy gradient theorem

Marginalize R, push in gradient

Dynamics + Reward constant w.r.t. $\theta$

Expanding $v_\pi(s')$ creates deeply nested computation:

At every step, compute every state you could get to from every state you could have been in

↓

Transform into simple sum over time steps and states:

What is total prob of being at each state at each time step?

$$\nabla v_\pi(s) = \nabla\left[\sum_a \pi(a|s)q_\pi(s,a)\right], \quad \text{for all } s \in \mathcal{S} \qquad \text{(Exercise 3.18)}$$

$$= \sum_a \left[\nabla\pi(a|s)q_\pi(s,a) + \pi(a|s)\nabla q_\pi(s,a)\right] \quad \text{(product rule of calculus)}$$

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

$$\text{(Exercise 3.19 and Equation 3.2)}$$

$$= \sum_a \left[\nabla\pi(a|s)q_\pi(s,a) + \pi(a|s)\sum_{s'} p(s'|s,a)\nabla v_\pi(s')\right] \qquad \text{(Eq. 3.4)} \longrightarrow$$

$$= \sum_a \left[\nabla\pi(a|s)q_\pi(s,a) + \pi(a|s)\sum_{s'} p(s'|s,a) \qquad \text{(unrolling)}\right.$$

$$\left.\sum_{a'}\left[\nabla\pi(a'|s')q_\pi(s',a') + \pi(a'|s')\sum_{s''} p(s''|s',a')\nabla v_\pi(s'')\right]\right]$$

$$= \sum_{x\in\mathcal{S}}\sum_{k=0}^{\infty} \Pr(s\to x, k, \pi)\sum_a \nabla\pi(a|x)q_\pi(x,a), \qquad \dashrightarrow$$

after repeated unrolling, where $\Pr(s\to x, k, \pi)$ is the probability of transitioning from state $s$ to state $x$ in $k$ steps under policy $\pi$. It is then immediate that

$$\nabla J(\boldsymbol{\theta}) = \nabla v_\pi(s_0)$$

$$= \sum_s \left(\sum_{k=0}^{\infty} \Pr(s_0\to s, k, \pi)\right)\sum_a \nabla\pi(a|s)q_\pi(s,a)$$

$$= \sum_s \eta(s)\sum_a \nabla\pi(a|s)q_\pi(s,a) \qquad \text{(box page 199)}$$

$$= \sum_{s'} \eta(s')\sum_s \frac{\eta(s)}{\sum_{s'}\eta(s')}\sum_a \nabla\pi(a|s)q_\pi(s,a)$$

$$= \sum_{s'} \eta(s')\sum_s \mu(s)\sum_a \nabla\pi(a|s)q_\pi(s,a) \qquad \text{(Eq. 9.3)}$$

$$\propto \sum_s \mu(s)\sum_a \nabla\pi(a|s)q_\pi(s,a) \qquad \text{(Q.E.D.)}$$

Handwritten annotations:

Proof of policy gradient theorem

Marginalize R, push in gradient

Dynamics + Reward constant w.r.t. $\theta$

Expanding $v_\pi(s')$ creates deeply nested computation:

At every step, compute every state you could get to from every state you could have been in

↓

Transform into simple sum over time steps and states:

What is total prob of being at each state at each time step?

unnormalized Steady-state prob of $s$ → $\eta(s)$

normalized version → $\mu(s)$

# REINFORCE

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

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

$$\boldsymbol{\theta}_{t+1} \doteq \boldsymbol{\theta}_t + \alpha \sum_a \hat{q}(S_t,a,\mathbf{w}) \nabla \pi(a|S_t,\boldsymbol{\theta})$$

All actions $\longrightarrow$   Q approx, not a sample return

# REINFORCE

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

$$\boldsymbol{\theta}_{t+1} \doteq \boldsymbol{\theta}_t + \alpha \sum_a \hat{q}(S_t, a, \mathbf{w}) \nabla \pi(a|S_t, \boldsymbol{\theta})$$

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

$$\nabla J(\boldsymbol{\theta}) = \mathbb{E}_\pi \left[ \sum_a \pi(a|S_t, \boldsymbol{\theta}) q_\pi(S_t, a) \frac{\nabla \pi(a|S_t, \boldsymbol{\theta})}{\pi(a|S_t, \boldsymbol{\theta})} \right]$$

$$= \mathbb{E}_\pi \left[ q_\pi(S_t, A_t) \frac{\nabla \pi(A_t|S_t, \boldsymbol{\theta})}{\pi(A_t|S_t, \boldsymbol{\theta})} \right] \qquad \text{(replacing } a \text{ by the sample } A_t \sim \pi)$$

$$= \mathbb{E}_\pi \left[ G_t \frac{\nabla \pi(A_t|S_t, \boldsymbol{\theta})}{\pi(A_t|S_t, \boldsymbol{\theta})} \right], \qquad \text{(because } \mathbb{E}_\pi[G_t|S_t, A_t] = q_\pi(S_t, A_t))$$

# REINFORCE

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

$$\boldsymbol{\theta}_{t+1} \doteq \boldsymbol{\theta}_t + \alpha \sum_a \hat{q}(S_t, a, \mathbf{w}) \nabla \pi(a|S_t, \boldsymbol{\theta})$$

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

$$\nabla J(\boldsymbol{\theta}) = \mathbb{E}_\pi \left[ \sum_a \pi(a|S_t, \boldsymbol{\theta}) q_\pi(S_t, a) \frac{\nabla \pi(a|S_t, \boldsymbol{\theta})}{\pi(a|S_t, \boldsymbol{\theta})} \right]$$

$$= \mathbb{E}_\pi \left[ q_\pi(S_t, A_t) \frac{\nabla \pi(A_t|S_t, \boldsymbol{\theta})}{\pi(A_t|S_t, \boldsymbol{\theta})} \right] \qquad \text{(replacing } a \text{ by the sample } A_t \sim \pi)$$

$$= \mathbb{E}_\pi \left[ G_t \frac{\nabla \pi(A_t|S_t, \boldsymbol{\theta})}{\pi(A_t|S_t, \boldsymbol{\theta})} \right], \qquad \text{(because } \mathbb{E}_\pi[G_t|S_t, A_t] = q_\pi(S_t, A_t))$$

---

**REINFORCE: Monte-Carlo Policy-Gradient Control (episodic) for $\pi_*$**

Input: a differentiable policy parameterization $\pi(a|s, \boldsymbol{\theta})$
Algorithm parameter: step size $\alpha > 0$
Initialize policy parameter $\boldsymbol{\theta} \in \mathbb{R}^{d'}$ (e.g., to $\mathbf{0}$)

Loop forever (for each episode):
    Generate an episode $S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot|\cdot, \boldsymbol{\theta})$
    Loop for each step of the episode $t = 0, 1, \ldots, T-1$:
        $G \leftarrow \sum_{k=t+1}^{T} \gamma^{k-t-1} R_k$            $(G_t)$
        $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \gamma^t G \nabla \ln \pi(A_t|S_t, \boldsymbol{\theta})$

# Gradient Bandits
### + Baseline

$$H_{t+1}(a) \doteq H_t(a) + \alpha \frac{\partial \mathbb{E}[R_t]}{\partial H_t(a)}$$

$$\frac{\partial \mathbb{E}[R_t]}{\partial H_t(a)} = \frac{\partial}{\partial H_t(a)} \left[ \sum_x \pi_t(x) q_*(x) \right]$$

$$= \sum_x q_*(x) \frac{\partial \pi_t(x)}{\partial H_t(a)}$$

$$= \sum_x \left( q_*(x) - B_t \right) \frac{\partial \pi_t(x)}{\partial H_t(a)},$$

Mean of Samples

Expectation Zero

# Gradient Bandits + Baseline

$$H_{t+1}(a) \doteq H_t(a) + \alpha \frac{\partial \mathbb{E}[R_t]}{\partial H_t(a)}$$

$$\frac{\partial \mathbb{E}[R_t]}{\partial H_t(a)} = \frac{\partial}{\partial H_t(a)} \left[ \sum_x \pi_t(x) q_*(x) \right]$$

$$= \sum_x q_*(x) \frac{\partial \pi_t(x)}{\partial H_t(a)}$$

$$= \sum_x \left( q_*(x) - B_t \right) \frac{\partial \pi_t(x)}{\partial H_t(a)},$$

Mean of Samples

Expectation Zero

# REINFORCE + Baseline

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

$$\nabla J(\boldsymbol{\theta}) \propto \sum_s \mu(s) \sum_a \left( q_\pi(s, a) - b(s) \right) \nabla \pi(a|s, \boldsymbol{\theta}).$$

$$\sum_a b(s) \nabla \pi(a|s, \boldsymbol{\theta}) = b(s) \nabla \sum_a \pi(a|s, \boldsymbol{\theta}) = b(s) \nabla 1 = 0.$$

$$\boldsymbol{\theta}_{t+1} \doteq \boldsymbol{\theta}_t + \alpha \left( G_t - \boxed{b(S_t)} \right) \frac{\nabla \pi(A_t|S_t, \boldsymbol{\theta}_t)}{\pi(A_t|S_t, \boldsymbol{\theta}_t)}.$$

$$\hat{v}(S_t)$$

**REINFORCE with Baseline (episodic), for estimating $\pi_\theta \approx \pi_*$**

Input: a differentiable policy parameterization $\pi(a|s,\boldsymbol{\theta})$
Input: a differentiable state-value function parameterization $\hat{v}(s,\mathbf{w})$
Algorithm parameters: step sizes $\alpha^{\boldsymbol{\theta}} > 0$, $\alpha^{\mathbf{w}} > 0$
Initialize policy parameter $\boldsymbol{\theta} \in \mathbb{R}^{d'}$ and state-value weights $\mathbf{w} \in \mathbb{R}^d$ (e.g., to $\mathbf{0}$)

Loop forever (for each episode):
    Generate an episode $S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot|\cdot,\boldsymbol{\theta})$
    Loop for each step of the episode $t = 0, 1, \ldots, T-1$:
        $G \leftarrow \sum_{k=t+1}^{T} \gamma^{k-t-1} R_k$                                       $(G_t)$
        $\delta \leftarrow G - \hat{v}(S_t,\mathbf{w})$
        $\mathbf{w} \leftarrow \mathbf{w} + \alpha^{\mathbf{w}} \delta \nabla \hat{v}(S_t,\mathbf{w})$
        $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha^{\boldsymbol{\theta}} \gamma^t \delta \nabla \ln \pi(A_t|S_t,\boldsymbol{\theta})$

---

**One-step Actor–Critic (episodic), for estimating $\pi_\theta \approx \pi_*$**

Input: a differentiable policy parameterization $\pi(a|s,\boldsymbol{\theta})$
Input: a differentiable state-value function parameterization $\hat{v}(s,\mathbf{w})$
Parameters: step sizes $\alpha^{\boldsymbol{\theta}} > 0$, $\alpha^{\mathbf{w}} > 0$
Initialize policy parameter $\boldsymbol{\theta} \in \mathbb{R}^{d'}$ and state-value weights $\mathbf{w} \in \mathbb{R}^d$ (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):
        $A \sim \pi(\cdot|S,\boldsymbol{\theta})$
        Take action $A$, observe $S', R$
        $\delta \leftarrow R + \gamma \hat{v}(S',\mathbf{w}) - \hat{v}(S,\mathbf{w})$         (if $S'$ is terminal, then $\hat{v}(S',\mathbf{w}) \doteq 0$)
        $\mathbf{w} \leftarrow \mathbf{w} + \alpha^{\mathbf{w}} \delta \nabla \hat{v}(S,\mathbf{w})$
        $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha^{\boldsymbol{\theta}} I \delta \nabla \ln \pi(A|S,\boldsymbol{\theta})$
        $I \leftarrow \gamma I$
        $S \leftarrow S'$

# Actor Only

- policy search
- Directly parameterized policy
- No value functions (except baseline in REINFORCE)

- Continuous actions natural to represent

- High variance, No bootstrapping

- Scales w/ policy complexity, not size of state space

## Actor Only

- policy search
- Directly parameterized policy
- No value functions (except baseline in REINFORCE)
- Continuous actions natural to represent
- High variance, No bootstrapping
- Scales w/ policy complexity, not size of state space

## Critic only

- value function methods
- Indirect policy via VF
- Discrete actions only
- Lower variance, bootstrapping
- Scales with size of state space

## Actor Only

- policy search
- Directly parameterized policy
- No value functions (except baseline in REINFORCE)
- Continuous actions natural to represent
- High variance, No bootstrapping
- Scales w/ policy complexity, not size of state space

## Critic only

- value function methods
- Indirect policy via VF
- Discrete actions only
- Lower variance, bootstrapping
- Scales with size of state space

## Actor - Critic

- Policy Search + value function
- Benefits of both!
- Continuous actions
- Bootstrapping
- Scales primarily with policy complexity

## Actor Only

- policy search
- Directly parameterized policy
- No value functions (except baseline in REINFORCE)
- Continuous actions natural to represent
- High variance, No bootstrapping
- Scales w/ policy complexity, not size of state space

## Critic only

- value function methods
- Indirect policy via VF
- Discrete actions only
- Lower variance, bootstrapping
- Scales with size of state space

## Actor - Critic

- Policy Search + value function
- Benefits of both!
- Continuous actions
- Bootstrapping
- Scales primarily with policy complexity

Many of most popular contempory methods are A-C:

- Proximal Policy Optimization
- A3C
- Soft Actor Critic
- DDPG
  ⋮