# Policy Gradient Methods

# Parameterized policy

- Learn a parameterized policy (must be differentiable) for action selection $\pi(a|s,\boldsymbol{\theta}) = \Pr\{A_t = a \mid S_t = s, \boldsymbol{\theta}_t = \boldsymbol{\theta}\}$

- A value function may still be used to learn the policy parameters (actor-critic)

- Optimisation by gradient accent in the direction of the gradient of $J(\theta)$ (some performance measure) with respect to $\theta$

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \alpha \widehat{\nabla J(\boldsymbol{\theta}_t)},$$

# Policy Approximation Advantages

- Parameterized numerical preferences $h(s, a, \theta)$ (i.e. neural network) with soft-max in action preferer

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

- Action preferences drive towards the optimal stochastic policy (and can approach a deterministic policy)
  - Action-value methods with $\varepsilon$-greedy action selection cannot
- Learn appropriate levels of exploration
- Policy may be a simpler function to approximate compared to Q-function
- Can input prior knowledge into policy

# Policy Gradient Theorem

- Episodic case, for which we define the performance measure as the value of the start state of the episode. $J(\boldsymbol{\theta}) \doteq v_{\pi_{\boldsymbol{\theta}}}(s_0),$

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

- Expression for how performance is affected by the policy parameter and doesn't involve derivatives of state distribution

- To convert to an algorithm all that is needed is some way of sampling whose expectation approximates this expression.

# REINFORCE: Monte Carlo Policy Gradient

- Replace sum over all states and actions by following target policy $\pi$ and sampling. Note, a weighting is introduced for an expectation under $\pi$

- Monto Carlo sampling so return $G_t$ from a complete episode

- REINFORCE update (stochastic gradient ascent)

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

- Update to increase the probability of taking action $A_t$ on future visits to $S_t$ proportional to the return and inversely proportional to the action probability

- Convergence to local optimate (decreasing $\alpha$) but high variance (slow learning)

**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 $t = 0, 1, \ldots, T-1$:

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

        $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \gamma^t G \nabla \ln \pi(A_t|S_t, \boldsymbol{\theta})$

# REINFORCE with Baseline

- Policy gradient theorem with comparison of action-value to a baseline (doesn't affect the gradient as long as $b(s)$ doesn't vary with $a$)

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

- REINFORCE with baseline update
$$\boldsymbol{\theta}_{t+1} \doteq \boldsymbol{\theta}_t + \alpha \Big(G_t - b(S_t)\Big) \frac{\nabla \pi(A_t|S_t, \boldsymbol{\theta}_t)}{\pi(A_t|S_t, \boldsymbol{\theta}_t)}.$$

- Baseline: state value estimate $\hat{v}(S_t, \mathbf{w})$,

- Reduces variance to speed up learning without introducing a bias (so will converge asymptotically to a local minimum)

**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 $t = 0, 1, \ldots, T-1$:

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

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

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

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

# Actor–Critic Methods

- Learn a value function to replace the full return with the one-step return (TD methods)

- State-value function assigns credit to the policy's action selections – critic the actor

- Bias introduced through bootstrapping (updating the value estimate for a state from the estimated values of subsequent states) reduces variance and speeds up learning

- Generalizations to n-step methods and then to a $\lambda$-return algorithm

**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–Critic with Eligibility Traces (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: trace-decay rates $\lambda^{\boldsymbol{\theta}} \in [0, 1]$, $\lambda^{\mathbf{w}} \in [0, 1]$; 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):
> $\quad$ Initialize $S$ (first state of episode)
> $\quad$ $\mathbf{z}^{\boldsymbol{\theta}} \leftarrow \mathbf{0}$ ($d'$-component eligibility trace vector)
> $\quad$ $\mathbf{z}^{\mathbf{w}} \leftarrow \mathbf{0}$ ($d$-component eligibility trace vector)
> $\quad$ $I \leftarrow 1$
> $\quad$ Loop while $S$ is not terminal (for each time step):
> $\quad\quad$ $A \sim \pi(\cdot|S, \boldsymbol{\theta})$
> $\quad\quad$ Take action $A$, observe $S', R$
> $\quad\quad$ $\delta \leftarrow R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})$ $\qquad$ (if $S'$ is terminal, then $\hat{v}(S', \mathbf{w}) \doteq 0$)
> $\quad\quad$ $\mathbf{z}^{\mathbf{w}} \leftarrow \gamma \lambda^{\mathbf{w}} \mathbf{z}^{\mathbf{w}} + \nabla \hat{v}(S, \mathbf{w})$
> $\quad\quad$ $\mathbf{z}^{\boldsymbol{\theta}} \leftarrow \gamma \lambda^{\boldsymbol{\theta}} \mathbf{z}^{\boldsymbol{\theta}} + I \nabla \ln \pi(A|S, \boldsymbol{\theta})$
> $\quad\quad$ $\mathbf{w} \leftarrow \mathbf{w} + \alpha^{\mathbf{w}} \delta \mathbf{z}^{\mathbf{w}}$
> $\quad\quad$ $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha^{\boldsymbol{\theta}} \delta \mathbf{z}^{\boldsymbol{\theta}}$
> $\quad\quad$ $I \leftarrow \gamma I$
> $\quad\quad$ $S \leftarrow S'$

# Policy Gradient for Continuing Problems

- For continuing problems we need to define performance in terms of the average rate of reward per time step:

$$J(\boldsymbol{\theta}) \doteq r(\pi) \doteq \lim_{h \to \infty} \frac{1}{h} \sum_{t=1}^{h} \mathbb{E}[R_t \mid S_0, A_{0:t-1} \sim \pi]$$

$$= \lim_{t \to \infty} \mathbb{E}[R_t \mid S_0, A_{0:t-1} \sim \pi]$$

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

- where $\mu$ is the steady-state distribution under $\pi$

- Value function are defined with respect to the differential return

$$G_t \doteq R_{t+1} - r(\pi) + R_{t+2} - r(\pi) + R_{t+3} - r(\pi) + \cdots.$$

# Policy Parameterization for Continuous Actions

- The policy can be defined as the normal probability density over a real-valued scalar action, with mean and standard deviation given by parametric function approximators that depend on the state.

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

- Parameter vector $\mu(s,\boldsymbol{\theta}) \doteq \boldsymbol{\theta}_\mu^\top \mathbf{x}_\mu(s)$ and $\sigma(s,\boldsymbol{\theta}) \doteq \exp\left(\boldsymbol{\theta}_\sigma^\top \mathbf{x}_\sigma(s)\right),$