# Policy Gradient

# Gole:

Given policy $\pi$ find best $\theta$.

$$\pi_\theta \longrightarrow \pi_{\theta*}$$

# How do we know which

# Policy is good

**or How to measure the quality of a policy $\pi_\theta$?**

## 1. Start state value:

$$J_1(\theta) = V^{\pi_\theta}(s_1) = \mathbb{E}[v_1]$$

## 2. Average value:

$$J_{avV}(\theta) = \sum_s d^{\pi_\theta}(s) V^{\pi_\theta}(s)$$

or

$$J_{avR}(\theta) = \sum_s d^{\pi_\theta}(s) \sum_a \pi_\theta(s, a) \mathcal{R}_s^a$$

# Policy Optimization

- Policy Based RL is an **optimization problem**.
- Find $\theta$ that maximises $J(\theta)$.

| Without Gradient | Uses Gradient |
| --- | --- |
| Hill climbing. | Gradient decent |
| Simplex /amoeba/ Nelder Mead. | Conjugate gradient |
| Genetic algorithms. | Quasi-newton |

- We will focus on gradient descent.

# Policy Gradient

Let $J(\theta)$ be andy policy objective function.

$\theta \rightarrow \theta'$ such that $J(\theta) < J(\theta')$

$$\Delta\theta = \alpha \nabla_\theta J(\theta)$$

Where

$$\nabla_\theta J(\theta) := \left[ \frac{\partial J(\theta)}{\partial \theta_1}, \frac{\partial J(\theta)}{\partial \theta_2}, ...., \frac{\partial J(\theta)}{\partial \theta_n} \right]$$

# Computing Gradients By Finite Differences

- To evaluate policy gradient of $\pi_\theta(s, a)$

- For each dimension $k \in [1, n]$

  - Estimate kth partial derivative of objective function w.r.t. $\theta$

  - perturbing $\theta$ by small amount $\epsilon$ in $k$ th dimension

$$\frac{\partial J(\theta)}{\partial \theta_k} \approx \frac{J(\theta + \epsilon u_k) - J(\theta)}{\epsilon}$$

- where $u_k$ is unit vector with 1 in kth component, 0 elsewhere

- Uses $n$ evaluations to compute policy gradient in $n$ dimensions

- Simple, noisy, inefficient - but sometimes effective

- Works for arbitrary policies, even if policy is not differentiable

7

# Score Function

$$\nabla_\theta \pi_\theta(s, a) = \pi_\theta(s, a) \ \frac{\nabla_\theta \pi_\theta(s, a)}{\pi_\theta(s, a)}$$

$$= \pi_\theta(s, a) \ \underbrace{\nabla_\theta \log_e(\pi_\theta(s, a))}_{\text{Score Function}}$$

# Softmax Policy

- We will use a softmax policy as a running example

- Weight actions using linear combination of features $\phi(s, a)^\top \theta$

- Probability of action is proportional to exponentiated weight

$$\pi_\theta(s, a) \propto e^{\phi(s,a)^\top \theta}$$

- The score function is

$$\nabla_\theta \log \pi_\theta(s, a) = \phi(s, a) - \mathbb{E}_{\pi_\theta}[\phi(s, \cdot)]$$

# Gaussian Policy

- In continuous action spaces, a Gaussian policy is natural
- Mean is a linear combination of state features $\mu(s) = \phi(s)^\top \theta$
- Variance may be fixed $\sigma^2$, or can also parametrised
- Policy is Gaussian, $a \sim \mathcal{N}\left(\mu(s), \sigma^2\right)$
- The score function is

$$\nabla_\theta \log \pi_\theta(s, a) = \frac{(a - \mu(s))\phi(s)}{\sigma^2}$$

# Monte-Carlo Policy Gradient (REINFORCE)

- Update parameters by stochastic gradient ascent
- Using policy gradient theorem
- Using return $v_t$ as an unbiased sample of $Q^{\pi_\theta}(s_t, a_t)$

$$\Delta\theta_t = \alpha \nabla_\theta \log \pi_\theta(s_t, a_t) v_t$$

**function REINFORCE**
    Initialise $\theta$ arbitrarily
    **for** each episode $\{s_1, a_1, r_2, ..., s_{T-1}, a_{T-1}, r_T\} \sim \pi_\theta$ **do**
        **for** $t = 1$ to $T - 1$ **do**
            $\theta \leftarrow \theta + \alpha \nabla_\theta \log \pi_\theta(s_t, a_t) v_t$
        **end for**
    **end for**
    **return** $\theta$
**end function**

# reducing variance using a critic

- monte-carlo policy gradient still has high variance
- we use a critic to estimate the action-value function,

$$q_w(s, a) \approx q^{\pi_\theta}(s, a)$$

- actor-critic algorithms maintain two sets of parameters
  - **Critic** Updates action-value function parameters $w$
  - **Actor** Updates policy parameters $\theta$, in direction suggested by critic
- Actor-critic algorithms follow an approximate policy gradient

$$\nabla_\theta J(\theta) \approx \mathbb{E}_{\pi_\theta}\left[\nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)\right]$$
$$\Delta\theta = \alpha \nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)$$

# Action-Value Actor-Critic

- Simple actor-critic algorithm based on action-value critic
- Using linear value fn approx. $Q_w(s, a) = \phi(s, a)^\top w$

    Critic Updates $w$ by linear TD(0)

    Actor Updates $\theta$ by policy gradient

**function** $\mathrm{QAC}$
    Initialise $s$, $\theta$
    Sample $a \sim \pi_\theta$
    **for** each step **do**
        Sample reward $r = \mathcal{R}_s^a$; sample transition $s' \sim \mathcal{P}_{s,\cdot}^a$
        Sample action $a' \sim \pi_\theta(s', a')$
        $\delta = r + \gamma Q_w(s', a') - Q_w(s, a)$
        $\theta = \theta + \alpha \nabla_\theta \log \pi_\theta(s, a) Q_w(s, a)$
        $w \leftarrow w + \beta \delta \phi(s, a)$
        $a \leftarrow a', s \leftarrow s'$
    **end for**
**end function**

# Critics at Different Time-Scales

- Critic can estimate value function $V_\theta(s)$ from many targets at different time-scales From last lecture...
  - For MC, the target is the return $v_t$
  
  $$\Delta\theta = \alpha(v_t - V_\theta(s))\phi(s)$$
  
  - For TD(0), the target is the TD target $r + \gamma V(s')$
  
  $$\Delta\theta = \alpha(r + \gamma V(s') - V_\theta(s))\phi(s)$$
  
  - For forward-view TD($\lambda$), the target is the $\lambda$-return $v_t^\lambda$
  
  $$\Delta\theta = \alpha(v_t^\lambda - V_\theta(s))\phi(s)$$
  
  - For backward-view TD($\lambda$), we use eligibility traces
  
  $$\delta_t = r_{t+1} + \gamma V(s_{t+1}) - V(s_t)$$
  $$e_t = \gamma\lambda e_{t-1} + \phi(s_t)$$
  $$\Delta\theta = \alpha\delta_t e_t$$

# Policy Gradient with Eligibility Traces

- Just like forward-view TD($\lambda$), we can mix over time-scales

$$\Delta\theta = \alpha(v_t^\lambda - V_v(s_t))\nabla_\theta \log \pi_\theta(s_t, a_t)$$

- where $v_t^\lambda - V_v(s_t)$ is a biased estimate of advantage fn
- Like backward-view TD($\lambda$), we can also use eligibility traces
  - By equivalence with TD($\lambda$), substituting $\phi(s) = \nabla_\theta \log \pi_\theta(s, a)$

$$\delta = r_{t+1} + \gamma V_v(s_{t+1}) - V_v(s_t)$$
$$e_{t+1} = \lambda e_t + \nabla_\theta \log \pi_\theta(s, a)$$
$$\Delta\theta = \alpha\delta e_t$$

- This update can be applied online, to incomplete sequences

# Summary of Policy Gradient Algorithms

- The policy gradient has many equivalent forms

$$\nabla_\theta J(\theta) = \mathbb{E}_{\pi_\theta} \left[ \nabla_\theta \log \pi_\theta(s, a) \, v_t \right] \qquad \text{REINFORCE}$$

$$= \mathbb{E}_{\pi_\theta} \left[ \nabla_\theta \log \pi_\theta(s, a) \, Q^w(s, a) \right] \qquad \text{Q Actor-Critic}$$

$$= \mathbb{E}_{\pi_\theta} \left[ \nabla_\theta \log \pi_\theta(s, a) \, A^w(s, a) \right] \qquad \text{Advantage Actor-Critic}$$

$$= \mathbb{E}_{\pi_\theta} \left[ \nabla_\theta \log \pi_\theta(s, a) \, \delta \right] \qquad \text{TD Actor-Critic}$$

$$= \mathbb{E}_{\pi_\theta} \left[ \nabla_\theta \log \pi_\theta(s, a) \, \delta e \right] \qquad \text{TD}(\lambda) \text{ Actor-Critic}$$

$$G_\theta^{-1} \nabla_\theta J(\theta) = w \qquad \text{Natural Actor-Critic}$$

- Each leads a stochastic gradient ascent algorithm
- Critic uses policy evaluation (e.g. MC or TD learning) to estimate $Q^\pi(s, a)$, $A^\pi(s, a)$ or $V^\pi(s)$