# Policy Gradient Method.

## Deep - Q - Learning

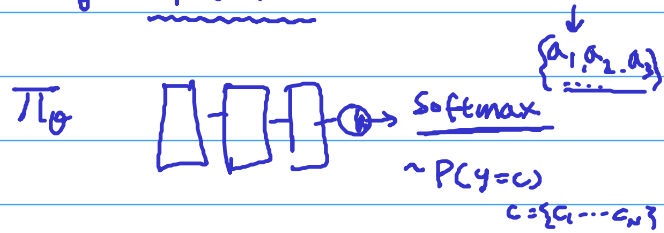$$\hat{Q}_\theta(s, a) = \underset{}{\cancel{\sum W_i f_i(s,a)}}$$

$$F(\theta, s, a)$$

$$\mathcal{L}(Q - \hat{Q}) =$$

$$\theta \leftarrow \theta - \alpha \frac{\partial \mathcal{L}}{\partial \theta}$$

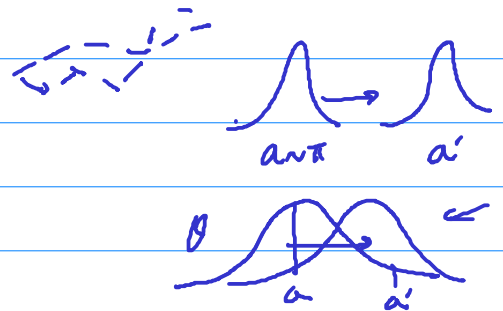Optimal policy $G_{\pi} = \sum_{t}^{T} \gamma^t \cdot r_t$

$$E[G_\theta]$$

$$Q \Rightarrow (s, a)$$

---

$$\pi_\theta = P(a|s)$$

$$\pi_\theta \quad \square - \square - \square \rightarrow \text{Softmax} \qquad \{a_1, a_2, a_3\}$$

$$\sim P(y=c)$$
$$c = \{c_1 \cdots c_N\}$$

$$P(a|s)$$

$\varepsilon$-greedy $\quad \pi \rightarrow$ random



$a \sim \pi \qquad a'$

$\theta$

$a \qquad a'$

PG: Stronger "convergence" guaranteed



$\mathcal{L}$

---

$$\pi_\theta = F_\theta(s)$$

$$\rightarrow J_\theta = \underset{\tau \sim \pi_\theta}{E}[r(\tau)], \quad \tau : \text{trajectory, rollout in an episode.}$$

$$\int P \cdot r(\tau) d\tau \qquad P(a|s)$$

$$= \int_\tau \pi_\theta(\tau) \cdot r(\tau) d\tau$$

$$\nabla_\theta J_\theta = \int_\tau \nabla_\theta \pi_\theta \, r(\tau) d\tau \qquad x \, d\log x = \frac{dx}{x}$$

$$= \int \pi_\theta \boxed{(\nabla_\theta \log \pi_\theta) \cdot r(\tau)} d\tau \qquad = \underset{\tau \sim \pi_\theta}{E}[\nabla_\theta \log \pi_\theta \cdot r(\tau)]$$

$$= \frac{1}{N} \sum_i^N \left( \sum_{t=0}^{T} \nabla_\theta \log \pi_\theta(a_{it}|s_{it}) \right) \cdot \left( \sum_{t=0}^{T} r(s_{it}, a_{it}) \right)$$

MC PG $\qquad \downarrow$

# REINFORCE



- gradient is noisy
- high variance

For $e$ in Sample $N$ episodes following $\pi_\theta$

$\quad$ For $t$ in $(0 \cdots T-1)$

$\quad\quad G_t \leftarrow \sum_{k=0}^{T} \gamma^{jk} R_k \leadsto t$

$\quad\quad \theta \leftarrow \theta + \alpha \nabla_\theta J \; (\cancel{\times})$

$J_\theta = \mathbb{E}_{\tau \sim \pi}[r(\tau)]$

$\boxed{\nabla_\theta J} = \underbrace{G \cdot \nabla_\theta \log \pi_\theta}_{}$

$\mathbb{E}_{\tau \sim \pi}[\underbrace{\nabla_\theta \log \pi_\theta) \cdot r(\tau)}_{g(\tau)}]$

$\boxed{\text{Var}(g(\tau) r(\tau)) = \mathbb{E}[g^2(\tau) r^2(\tau)] - \mathbb{E}[g \cdot r]^2}$

## 1) Causality

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

$\sum_{t=0}^{T} r = \sum_{k=0}^{t} r + \sum_{t}^{T} r \quad$ reward to go $\; (G_t)$

$\nabla_\theta J_\theta \sim G_t$

## 2) baseline

$J = \mathbb{E}_{\tau \sim \pi_\theta}[r(\tau)] \rightarrow \mathbb{E}_{\tau \sim \pi_\theta}[r(\tau) - b]$

$\boxed{\nabla_\theta J} \rightarrow$ Policy Gradient $\nabla_\theta \log \pi_\theta \quad$ or $\quad \nabla_\theta \pi_\theta$.

$\nabla_\theta J = \mathbb{E}_{\tau \sim \pi_\theta}[\nabla_\theta \log \pi_\theta \cdot (r - b)] = \mathbb{E}_{\tau \sim \pi_\theta}[\nabla_\theta \log \pi_\theta \cdot r]$

$\quad = \mathbb{E}_{\tau \sim \pi}[\nabla_\theta \log \pi_\theta \, r] - b \underbrace{\mathbb{E}_{\tau \sim \pi}[\nabla_\theta \log \pi_\theta]}_{\rightarrow \, \sim 0}$

$b \rightarrow \underbrace{\frac{1}{N} \sum_{i}^{N} r_i(\tau)}_{\text{// practical}}$

$$\text{Var}\left[\underbrace{\nabla_\theta \log \pi_\theta}_{g(\tau)} \underbrace{(r-b)}_{(r(\tau)-b)}\right]$$

$$= \underset{\tau\sim\pi}{E}\left[g^2(r-b)^2\right] - \underset{\tau\sim\pi}{E}\left[g\cdot(r-b)\right]^2 \qquad E[g(r-b)] = E[g\,r]$$

$$\text{Var} = E\left[g^2 r^2\right] - E\left[2\cdot g\cdot r\cdot b\right] - E\left[g\cdot r\right]^2$$
$$+ E\left[g^2 b^2\right]$$

$$b \qquad \frac{\partial \text{Var}}{\partial b} = 0 = -\frac{\partial}{\partial b}E\left[2\cdot g\cdot r\cdot b\right] + \frac{\partial}{\partial b}E\left[g^2 b^2\right] = 0$$

$$-E\left[2\cdot g\cdot r\right] + 2b\,E\left[g^2\right] = 0$$

$$b = \frac{E\left[g(\tau)^2\cdot r(\tau)\right]}{E\left[g(\tau)^2\right]} \nearrow$$