

$$s \in S, a \in A, \pi(a|s), p(s'|s, a), r(s, a)$$

$$\tau = \{s_0, a_0, r_0, s_1, a_1, r_1, \dots\}$$

$$p(\tau) = p(s_0) \pi(a_0|s_0) p(s_1|s_0, a_0) \pi(a_1|s_1) \dots$$

$$R_0 = r_0 + \gamma r_1 + \gamma^2 r_2 + \dots \quad \gamma \in (0, 1)$$

$$\mathbb{E}_{p(\tau|\pi)} R_0 \rightarrow \max_{\pi}$$

$$V^{\pi}(s) = \mathbb{E}[R_0 \mid s_0 = s]$$

$$Q^{\pi}(s, a) = \mathbb{E}[R_0 \mid s_0 = s, a_0 = a]$$

$$V^{\star}(s) = \max_{\pi} V^{\pi}(s)$$

$$Q^{\star}(s, a) = \max_{\pi} Q^{\pi}(s, a)$$

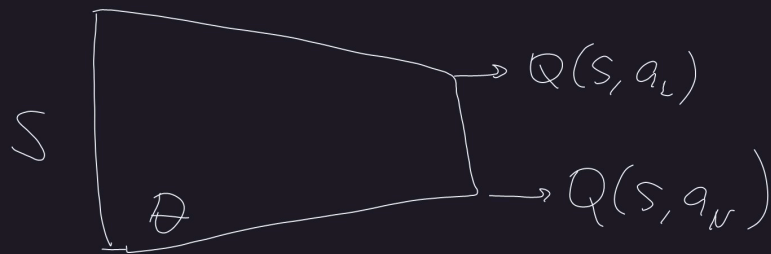
$$\forall s, a \quad Q^{\star}(s, a) = r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} \max_{a'} Q^{\star}(s', a')$$

$$\pi(a|s) = \arg \max_a Q^{\star}(s, a)$$

$$Q(s, a | \theta)$$

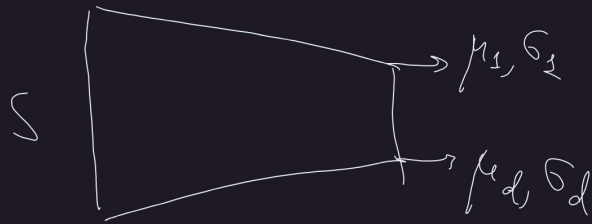
$$\text{Loss}(\theta) = \mathbb{E}_{s, a, s', r} \left(Q(s, a | \theta) - \left(r + \gamma \max_{a'} Q(s', a' | \bar{\theta}) \right) \right)^2$$

$\rightarrow \min_{\theta}$



Policy Gradient RL

$$\pi(a|s, \theta) = \prod_{j=1}^d \mathcal{N}(a_j | \mu_j(s, \theta), \log(1 + \exp(\sigma_j(s, \theta))))$$



$$\pi(a|s, \theta) = \text{Softmax}(\text{outputs}(s, \theta))$$

$$J(\theta) = \mathbb{E}_{p(\tau|\theta)} R(\tau) \rightarrow \max_{\theta}$$

Log-deriv. trick:

$$\nabla_{\theta} J(\theta) =$$

$$= \nabla_{\theta} \int p(\tau|\theta) R(\tau) d\tau =$$

$$= \int \nabla_{\theta} p(\tau|\theta) R(\tau) d\tau =$$

$$= \int p(\tau|\theta) \left(\frac{\nabla_{\theta} p(\tau|\theta)}{p(\tau|\theta)} \right) R(\tau) d\tau =$$

$$= \mathbb{E}_{p(\tau|\theta)} \nabla_{\theta} \log p(\tau|\theta) \cdot R(\tau)$$

Reparameterization trick:

$$\nabla_{\mu} \mathbb{E}_{N(x|\mu, \sigma^2)} f(x) =$$

$$= \nabla_{\mu} \mathbb{E}_{N(\varepsilon|0, 1)} f(\mu + \sigma \varepsilon) =$$

$$= \mathbb{E}_{N(\varepsilon|0, 1)} \nabla_{\mu} f(\mu + \sigma \varepsilon)$$

$$\tau = \{ s_0, a_0, z_0, \dots, s_{T-1}, a_{T-1}, z_{T-1}, s_T \}$$

$$\mathbb{E}_{p(\tau|\theta)} \nabla_{\theta} \log p(\tau|\theta) \left[\sum_{t=0}^T \gamma^t z_t \right] =$$

$$= \left\{ p(\tau|\theta) = p(s_0) \prod (a_u | s_u, \theta) p(s_{\underline{1}} | s_0, a_0) \prod (a_u | s_u, \theta) \dots \right\} =$$

$$= \mathbb{E}_{p(\tau|\theta)} \left[\sum_{u=0}^T \nabla \log \prod (a_u | s_u, \theta) \right] \left[\sum_{t=0}^T \gamma^t z_t \right] =$$

$$= \sum_{u,t} \mathbb{E}_{p(\tau|\theta)} \nabla_{\theta} \log \prod (a_u | s_u, \theta) \cdot \gamma^t z_t \quad u > t$$

$$\underline{u > t}$$

$$\begin{aligned} & \mathbb{E}_{s_0, a_0, s_1, a_1, \dots} \nabla \log \pi(a_u | s_u, \theta) \gamma^t z_t = \\ &= \mathbb{E}_{s_0, a_0, \dots, \color{blue}{s}_{u-1}, a_{u-1}, s_u} \gamma^t z_t \underbrace{\mathbb{E}_{\pi(a_u | s_u, \theta)} \nabla \log \pi(a_u | s_u, \theta)}_{11} \end{aligned}$$

$$\gamma(s, a)$$

$$\begin{aligned} & \sum_{a_u} \pi(a_u | s_u, \theta) \nabla \log \pi(a_u | s_u, \theta) = \\ &= \sum_{a_u} \cancel{\pi(a_u | s_u, \theta)} \frac{\nabla \pi(a_u | s_u, \theta)}{\cancel{\pi(a_u | s_u, \theta)}} \overset{11}{=} \nabla_{\theta} \sum_{a_u} \pi(a_u | s_u, \theta) = \\ & \qquad \qquad \qquad = 0 \end{aligned}$$

$$\nabla_{\theta} J(\theta) = \sum_{t=0}^T \mathbb{E}_{p(\tau|\theta)} \nabla_{\theta} \log \pi(a_t | s_t, \theta) \cdot \left(\sum_{u=t}^T \gamma^u r_u \right)$$

↑
weighted ML maximization

REINFORCE algorithm

Init. θ

For episodes:

Init. s_0

For $t=0, 1, \dots, T$:

$a_t \sim \pi(a_t | s_t, \theta)$

get r_t, s_{t+1}

$R := 0, L_0 = 0$

For $t=T, T-1, \dots, 0$:

$R = r_t + \gamma R$

$L_0 = L_0 - \log \pi(a_t | s_t, \theta) \cdot \gamma^t \cdot R$

$\theta = \text{opt_step}(\nabla_{\theta} L)$

Limitations:

- full episodes
- huge variance of SG

$$J(\theta) = \mathbb{E}_{p(\tau|\theta)} R(\tau) = \mathbb{E}_{p(s_0)} V^{\pi_\theta}(s_0)$$

Actor-Critic

$\pi(a|s, \theta_\pi)$ - actor

$Q(s, a | \theta_Q)$ - critic

Th (PG theorem)

$$J(\theta) = \mathbb{E}_{p(s_0)} \mathbb{E}_{\pi(a_0|s_0, \theta)} Q^{\pi_\theta}(s_0, a_0)$$

$$\Rightarrow \nabla_{\theta} J(\theta) = \frac{1}{1-\gamma} \mathbb{E}_{p(s_0)} \mathbb{E}_{s \sim p^{\pi_\theta}} \mathbb{E}_{a \sim \pi(a|s, \theta)} \nabla_{\theta} \log \pi(a|s, \theta) Q^{\pi_\theta}(s, a)$$

$$E_{\pi(a|s, \theta)} \sum_{\theta} \log \pi(a|s, \theta) (Q(s, a) - \overset{\text{baseline}}{B(s)})$$

$$\hookrightarrow E_{\pi(a|s, \theta)} \sum_{\theta} \log \pi(a|s, \theta) (B(s)) = B(s) \cdot 0$$

Baseline choice $B(s) = E_{\pi(a|s, \theta)} Q^{\pi}(s, a) = V^{\pi}(s)$

Example $E_a f_1(a) \cdot f_2(a)$

$$a: \begin{array}{cc} 0 & 1 \\ \frac{1}{2} & \frac{1}{2} \end{array} \quad f_1 = \begin{cases} 1, & \text{if } a=0 \\ -1, & \text{if } a=1 \end{cases}$$

$$f_2 = \begin{cases} 100, & \text{if } a=0 \\ 101, & \text{if } a=1 \end{cases}$$

$$E_a f_1 = \frac{1}{2} \cdot 1 + \frac{1}{2} \cdot (-1) = 0$$

$$E_a f_1 \cdot f_2 = \frac{1}{2} (1 \cdot 100) + \frac{1}{2} ((-1) \cdot 101) = -\frac{1}{2}$$

$$E_a (f_1 \cdot f_2 - E f_1 f_2)^2 = E_a \left(f_1 \cdot f_2 + \frac{1}{2} \right)^2 =$$

$$= \frac{1}{2} \left(1 \cdot 100 + \frac{1}{2} \right)^2 + \frac{1}{2} \left((-1) \cdot 101 + \frac{1}{2} \right)^2 = (100.5)^2$$

$$E_a f_1 (f_2 - E f_2) = E f_1 f_2$$

$$E f_2 = \frac{1}{2} \cdot 100 + \frac{1}{2} \cdot 101 = 100.5$$

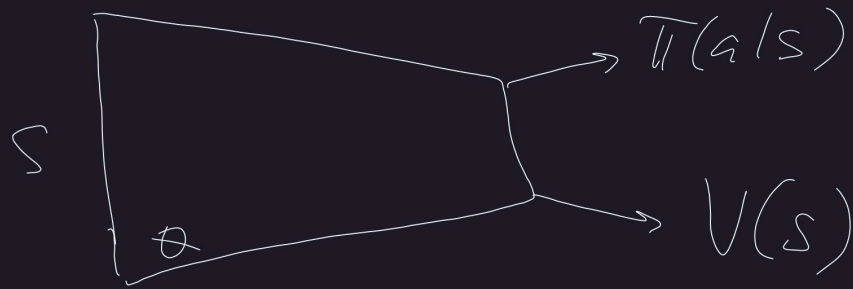
$$Var = E_a \left(f_1 (f_2 - 100.5) + \frac{1}{2} \right)^2 =$$

$$= \frac{1}{2} \left(1 \cdot (100 - 100.5) + \frac{1}{2} \right)^2 +$$

$$+ \frac{1}{2} \left((-1) \cdot (101 - 100.5) + \frac{1}{2} \right)^2 = 0$$

$$A^\pi(s, a) := Q^\pi(s, a) - V^\pi(s)$$

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)} V^\pi(s')$$



A2C (Advantage Actor Critic)

Init. θ

For episodes:

$k = 0$

while $k < T$:

for $t = k, k+1, \dots, k+n$:

$a_t \sim \pi(a_t | s_t, \theta)$

get r_t, s_{t+1}

$R = V(s_{k+n} | \bar{\theta}), L_{\theta_\pi} = 0, L_{\theta_v} = 0$

for $t = k+n-1, \dots, k$:

$R = r_t + \gamma \cdot R$

$L_v = L_v + \frac{1}{2} (V(s_t | \bar{\theta}) - R)^2$

$L_{\theta_\pi} = L_{\theta_\pi} - \log \pi(a_t | s_t, \theta) (R - V(s_t | \bar{\theta}))$

use parallel
agent



$\pi(a | s, \theta)$ - learned policy

$b(s)$ - behaviour policy

If $b(s) = \pi(a | s, \theta)$

\Rightarrow on-policy RL

Otherwise off-policy RL

$\theta_v = \text{opt_step}(\nabla L_v)$

$\theta_\pi = \text{opt_step}(\nabla L_{\theta_\pi})$



RLHF (RL from Human Feedback)

$$\tau_1 = \{s_0^1, a_0^1, s_1^1, a_1^1, \dots, s_n^1, a_n^1\}$$

$$\tau_2 = \{s_0^2, a_0^2, s_1^2, a_1^2, \dots, s_n^2, a_n^2\}$$

$$\tau_1 \cup \tau_2 \rightarrow \begin{cases} \mu = (1, 0) & \text{if } \tau_1 > \tau_2 \\ \mu = (0, 1) & \text{if } \tau_1 < \tau_2 \\ \mu = (\frac{1}{2}, \frac{1}{2}) & \text{if } \tau_1 \approx \tau_2 \end{cases}$$

$$r_{\theta}(s, a)$$

$$p(\hat{v}_1 > \hat{v}_2 | \theta) = 0.9 \frac{\exp\left(\sum_{s,a \in \hat{v}_1} r_{\theta}(s, a)\right)}{\exp\left(\sum_{s,a \in \hat{v}_1} r_{\theta}(s, a)\right) + \exp\left(\sum_{s,a \in \hat{v}_2} r_{\theta}(s, a)\right)} \sqrt{+ 0.1 \cdot \frac{1}{2}}$$

$$= \sigma\left(\sum_{s,a \in \hat{v}_1} r_{\theta}(s, a) - \sum_{s,a \in \hat{v}_2} r_{\theta}(s, a)\right)$$

$$\mathcal{L}(\theta) = -\mathbb{E}_{(\hat{v}_1, \hat{v}_2, \mu)} \left(\mu(1) \log p(\hat{v}_1 > \hat{v}_2 | \theta) + \mu(2) \log p(\hat{v}_2 > \hat{v}_1 | \theta) \right) \rightarrow \min_{\theta}$$