

Policy Gradient - II

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \nabla_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau) \\ &\approx \frac{1}{N} \sum_{i=1}^N \underbrace{R(\tau_i)}_{\text{blue box}} \underbrace{\nabla_{\theta} \log P(\tau_i; \theta)}_{\text{red box}} \quad - (1) \end{aligned}$$

where

$$R(\tau_i) = \sum_{t=0}^{\tau_i-1} r(s_t, a_t)$$

$$\nabla_{\theta} \log P(\tau_i; \theta) = \sum_{t=0}^{\tau_i-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

There is Causality violation issue

"Action in the future can not affect rewards in the past"

* Causality

For example, let's assume we have single reward, at j -th time step.

$$\nabla_{\theta} \mathbb{E}_{s_0, a_0, \dots, \underbrace{s_{\tau-1}}_{\text{future}}} [r_j] = \nabla_{\theta} \mathbb{E}_{s_0, a_0, \dots, s_j} [r_j]$$

the expectation stops at j -th term, all other terms cancel out.

Using the linearity of expectation, such as

$$\nabla_{\theta} \mathbb{E}_{\tau} \left[\sum_{t=0}^{\tau-1} r(s_t, a_t) \right] = \nabla_{\theta} \sum_{t=0}^{\tau-1} \mathbb{E}_{\tau} [r(s_t, a_t)]$$

so, only causal terms matter,

$$\nabla_{\theta} \mathbb{E}_{\pi} [r(s_j, a_j)] = \mathbb{E}_{s_0, a_0 \dots s_j, a_j} \left[r(s_j, a_j) \sum_{t=0}^j \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \right]$$

the reward is only affected by actions that came before,

$$g_i = \sum_{t'=0}^{T-1} r(s_{t'}, a_{t'}) \sum_{t=0}^{t'} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \quad \text{--- (2)}$$

$$g_i = \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \left(\sum_{t'=t}^{T-1} r_{t'} \right) \quad \text{--- (3)}$$

← Simplified version

↳ $r(s_{t'}, a_{t'})$

(2) → (3)

$$g_i = \sum_{t'=0}^{T-1} \underbrace{r(s_{t'}, a_{t'})}_{r_{t'}} \underbrace{\sum_{t=0}^{t'} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)}_{h_t}$$

$t'=0$	$r_0 h_0$		
$t'=1$	$r_1 h_0$	$+ r_1 h_1$	
$t'=2$	$r_2 h_0$	$+ r_2 h_1$	$+ r_2 h_2$
\vdots			
$t'=T-1$	$r_{T-1} h_0$	$+ r_{T-1} h_1$	\dots

$$\left(\sum_{t=0}^{T-1} r_t \right) h_0 + \left(\sum_{t=1}^{T-1} r_t \right) h_1 + \dots + \left(\sum \dots \right) h_{T-1}$$

$$\Leftrightarrow \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \left(\sum_{t'=t}^{T-1} r(s_{t'}, a_{t'}) \right)$$

Quick Review.

Causal

$$g_i = \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \left(\sum_{t'=t}^{T-1} r_{t'} \right)$$

$r(s_{t'}, a_{t'})$
↑

Non-Causal

$$g_i = \left(\sum_{t=0}^{T-1} r(s_t, a_t) \right) \left(\sum_{t'=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_{t'}, s_{t'}) \right)$$

Gradient Update

$$g \approx \frac{1}{N} \sum_{i=1}^N g_i$$

$$\theta_{k+1} \leftarrow \theta_k + \alpha g.$$

REINFORCE algorithm (1992, Ronald Williams)

- Initialize policy θ_0 , and learning rate α
- For $i=1 : \text{num-Iter}$:

For $j=1 : \text{num-rolls}$:

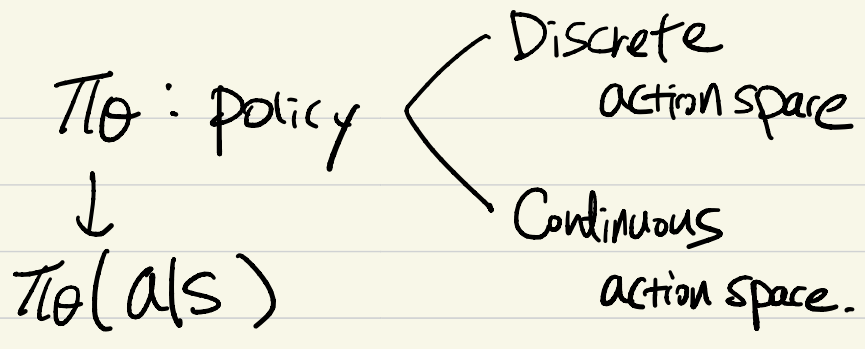
Compute the grad. estimate $g_i = \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \left(\sum_{t'=t}^{T-1} r_{t'} \right)$

Estimate $g \approx \frac{1}{N} \sum_{i=1}^N g_i$

Gradient update: $\theta_{k+1} \leftarrow \theta_k + \alpha g.$

What's

" $\nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$ " ?



* Diagonal Gaussian Policy.

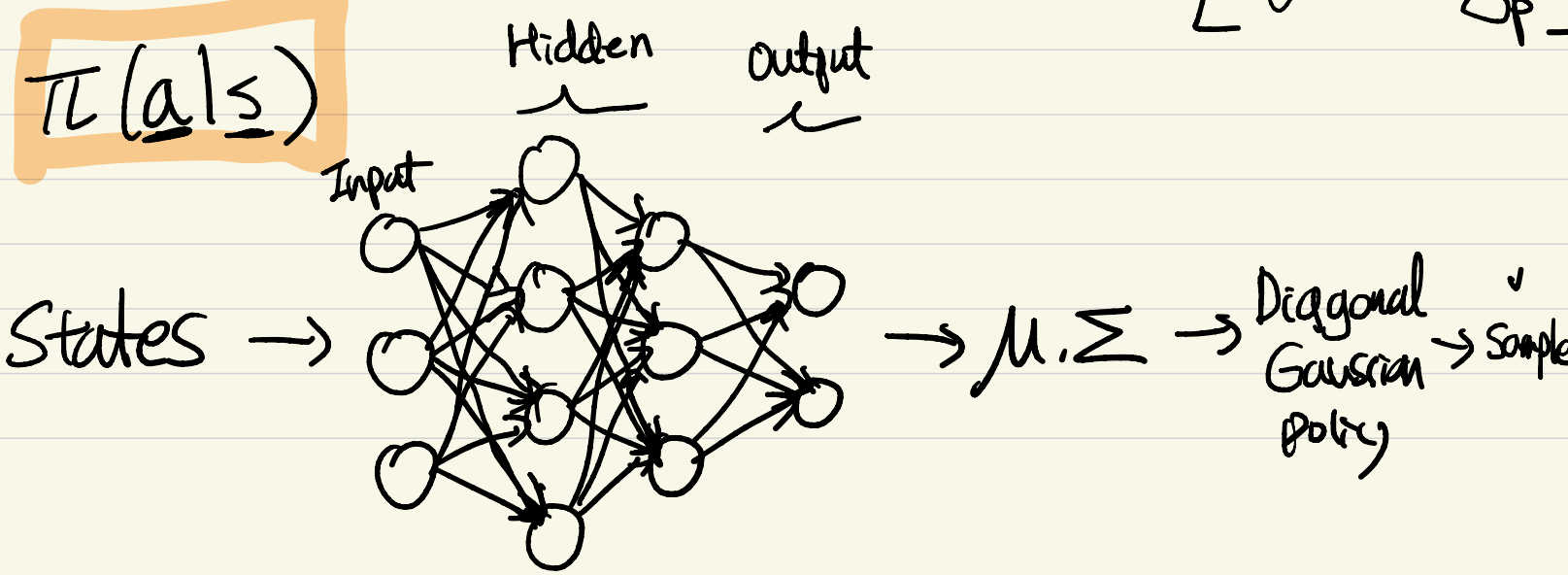
↳ p-dimensional Gaussian policy, $|A| = p$

NN takes in states and outputs mean, $\mu \in \mathbb{R}^{|A|}$ and covariance, $\Sigma \in \mathbb{R}^{|A|}$

$$\pi(a|\mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^p \det(\Sigma)}} \exp^{-\frac{1}{2}(a-\mu)^T \Sigma^{-1} (a-\mu)}$$

where Σ is covariance matrix, $\Sigma = \begin{bmatrix} \sigma_1^2 & & 0 \\ & \sigma_2^2 & \\ 0 & & \ddots \\ & & & \sigma_p^2 \end{bmatrix}$

$\pi(a|s)$



$\nabla_{\theta} \log \pi_{\theta}(a_t | S_t)$:

$$\bullet \pi(a | \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^P \det(\Sigma)}} \exp^{-\frac{1}{2}(a-\mu)^T \Sigma^{-1} (a-\mu)}$$

$$\bullet \log \pi_{\theta}(a | \mu, \Sigma) = -\frac{1}{2} \left(P \log(2\pi) + \log \left(\prod_i \sigma_i^2 \right) \right) - \frac{1}{2} (a-\mu)^T \Sigma^{-1} (a-\mu)$$

$$\bullet \nabla_{\theta} \log \pi_{\theta}(a | \mu, \Sigma) = -\frac{1}{2} \left(\nabla_{\theta} \log \left(\prod_i \sigma_i^2 \right) \right) - \nabla_{\theta} \frac{1}{2} (a-\mu)^T \Sigma^{-1} (a-\mu)$$

↳ Automatic differentiation tool, i.e., tensorflow
Computes automatically.

