

# RL 2022 - lec3 : policy gradient

The setup :

$$\max_{\rho} J(x_0) = \mathbb{E} \left[ \sum_{k=0}^N \gamma^k r(X_k, V_k) \mid X_0 = x_0 \right]$$

$$\text{s.t. } X_{k+1} \sim f(x_{k+1} \mid x_k, u_k)$$

$$V_k \sim p^\theta(u_k \mid x_k)$$

$\theta$  — policy par. (tensor)

$$Z_k := (X_k, V_k)$$

$\Rightarrow$  can rewrite the env. as

$$Z_{k+1} \sim P^\theta(Z_{k+1} \mid Z_k)$$

State-action trajectory :

$$\bar{Z}_k^N := \{Z_k, Z_{k+1}, \dots, Z_{k+N-1}\}$$

The accumulated reward  
(return) :

$$G_k^N(\bar{Z}_k^N) := \sum_{i=k}^{k+N-1} \gamma^i r(Z_i)$$

Let's unwrap the objective :

$$J^\theta(x_0) = \mathbb{E} \left[ G_o^N(\bar{Z}_o^N) \mid X_o = x_0 \right]$$

$$= \int G_o^N(\bar{Z}_o^N) P[d\bar{Z}_o^N \mid X_o = x_0]$$

X<sup>N</sup> × U<sup>N</sup>
↑  
state space    action space    integration w.r.t.  
prob-measure

In finite state and action space case, that integral would amount to

$$\sum_j G_o^N(\bar{Z}_o^{N(j)}) P[\bar{Z}_o^{N(j)}]$$

index of an outcome

An example policy (Gaussian)

$$\textcircled{1} \quad p^{\theta}(u|x) = \frac{1}{\sigma(x, \theta)\sqrt{2\pi}} e^{-\frac{(u - \mu(x, \theta))^2}{2\sigma^2(x, \theta)}},$$

where feature vector

$$\mu(x, \theta) = \theta_u^T \varphi_u(x),$$

$$\sigma(x, \theta) = e^{\theta_\sigma^T \varphi_\sigma(x)}$$

feature vector

$$\theta = (\theta_u, \theta_\sigma)$$

\textcircled{2} In case of finite action

space, say,

soft-argmax

$$p^{\theta}(u|x) = \frac{e^{\kappa(u|x, \theta)}}{\sum_{s \in \mathcal{U}} e^{\kappa(s|x, \theta)}},$$

where  $\kappa(u|x, \theta) = \theta^T \varphi(u, x)$

Getting back:

$$J^\theta(x_0) = \int_{\mathcal{X}^N \times \mathcal{U}^N} G_o^N(\bar{z}_o^N) \underbrace{\mathbb{P}[d\bar{z}_o^N]}_{}$$

$$= \int_{\mathcal{X}^N \times \mathcal{U}^N} G_o^N(\bar{z}_o^N) P_o(x_0) \cdot \prod_{k=0}^N p^\theta(u_k|x_k) \cdot f(x_{k+1}|x_k, u_k) d\bar{z}_o^N$$

---

Let's use gradient ascent rule  
for the policy update:

$$\theta_{i+1} := \theta_i + \lambda \nabla_\theta J^\theta(x_0) \Big|_{\theta=\theta_i}$$

$\uparrow$                        $\uparrow$   
gradient step        learning rate

---

Log - likelihood trick:

$$\nabla_\theta P^\theta(z) = P^\theta(z) \nabla_\theta \ln P^\theta(z)$$

Now, apply the trick to our case :

$$\nabla_{\theta} \left( \int_{\mathcal{X}^N \times \mathcal{U}^N} G_o^N(\bar{z}_o^N) P_o(x_o) \cdot \prod_{k=0}^N p^{\theta}(u_k | x_k) \cdot f(x_{k+1} | x_k, u_k) d\bar{z}_o^N \right) =$$

$$\int_{\mathcal{X}^N \times \mathcal{U}^N} G_o^N(\bar{z}_o^N) P_o(x_o) \cdot \prod_{k=0}^N p^{\theta}(u_k | x_k) \cdot f(x_{k+1} | x_k, u_k)$$

$$\nabla_{\theta} \ln \left( P_o(x_o) \cdot \prod_{k=0}^N p^{\theta}(u_k | x_k) \cdot f(x_{k+1} | x_k, u_k) \right) d\bar{z}_o^N =$$

Observe the following :

$$\mathbb{E}[h(x)] = \int_{\mathcal{X}} P(x) h(x) dx = \int h(x) \mathbb{P}[dx]$$

$$\int_{\mathcal{X}^N \times \mathcal{U}^N} G_o^N(\bar{z}_o^N) P_o(x_o) \cdot \prod_{k=0}^N p^{\theta}(u_k | x_k) \cdot f(x_{k+1} | x_k, u_k)$$

$$\nabla_{\theta} \ln \left( P_o(x_o) \cdot \prod_{k=0}^N p^{\theta}(u_k | x_k) \cdot f(x_{k+1} | x_k, u_k) \right) d\bar{z}_o^N =$$

$$\mathbb{E} \left[ G_o^N(\bar{Z}_o^N) \sum_{k=0}^N \nabla_\theta \ln \rho^\theta(V_u | X_u) \right]$$

Plugging this into the gradient ascent rule gives REINFORCE:

$$\begin{aligned}\theta_{i+1} &:= \theta_i + \alpha \mathbb{E} \left[ G_o^N(\bar{Z}_o^N) \sum_{k=0}^N \nabla_\theta \ln \rho^\theta(V_u | X_u) \right] \\ &= \theta_i + \alpha \mathbb{E} \left[ \sum_{k=0}^N \nabla_\theta \ln \rho^\theta(V_u | X_u) \cdot \sum_{i=0}^N \delta^i \pi(x_i, V_i) \right]\end{aligned}$$