

## Lecture 25

CS59000 - RL

Model free policy learning

Agenda:

- policy gradient

---

Intuition:  $M$ , Fixed horizon

Given a policy  $\pi$

$$J(\pi) = E_{\pi} \left[ \sum_{h=1}^H r_h \right] \Rightarrow \Psi = \sum_{h=1}^H r_h$$

$$\tau = (x_1, a_1, r_1, \dots, x_H, a_H, r_H)$$

$$J(\pi) = \sum Pr(\tau; \pi) \Psi_{\tau} \quad ; \quad \Psi_{\tau} \Rightarrow \text{reward in trajectory } \tau.$$

Let assume reward is deterministic

$$\Rightarrow r(x, a)$$

$$\Rightarrow \tau = (x_1, a_1, x_2, a_2, \dots)$$

$$\Rightarrow \Psi(\tau) = \sum_{h=1}^H r(x_h, a_h)$$

$$\Rightarrow Pr(\tau; \pi) = P_1(x_1) \pi(a_1; x_1) P(x_2 | x_1, a_1) \\ \dots \pi(a_H; x_H)$$

$$\eta(\pi) = \sum_{\tau} P_{\tau}(\tau; \pi) \psi(\tau)$$

$$\Rightarrow \nabla_{\pi} \eta(\pi) = \nabla_{\pi} \left( \sum_{\tau} P_{\tau}(\tau; \pi) \psi(\tau) \right)$$

$$= \sum_{\tau} \nabla_{\pi} P_{\tau}(\tau; \pi) \psi(\tau)$$

$$\nabla \log \pi = \frac{\nabla \pi}{\pi}$$

$$= \sum_{\tau} P_{\tau}(\tau; \pi) \nabla_{\pi} \log P_{\tau}(\tau; \pi) \psi(\tau)$$

$$= E \left[ \nabla_{\pi} \log P_{\tau}(\tau; \pi) \psi(\tau) \right]$$

$$\ast \log P_{\tau}(\tau; \pi) = \underbrace{\log p_1(x_1)}_{\tau} + \underbrace{\log \pi(A_1; x_1)}_{\tau} + \underbrace{\log p(x_2 | x_1, A_1)}_{\tau}$$

$$\Rightarrow \nabla_{\pi} \log P_{\tau}(\tau; \pi) = \sum_{h=1}^H \log \pi(A_h; x_h)$$

$$\Rightarrow \nabla_{\pi} \eta(\pi) = E \left[ \nabla_{\pi} \left( \sum \log \pi(A_h; x_h) \right) \psi(\tau) \right]$$

Monte Carlo sampling

$$\Rightarrow \nabla_{\pi} \hat{\eta}(\pi) = \frac{1}{N} \sum_{i=1}^N \left( \nabla_{\pi} \left( \sum_{h=1}^H \log \pi(A_h^i; x_h^i) \right) \psi(\tau^i) \right)$$

Consider an MDP  $M: (X, A, P, P_0, R)$

- $X$  as a Borel space (state space)
- $A$  as a Borel space (Action space)
- $P_0$  a probability measure on the initial state.
- $\bar{P}$  a transition kernel

-  $R$  a reward kernel,  $\bar{r}(x, a)$  it means  $\forall x, y \in X \times A$

Consider a set of parameters  $\Theta$  such that each  $\theta \in \Theta$  parametrizes a policy kernel  $\bar{\pi}_\theta$  from

$(X, \mathcal{B}(X))$  to  $(A, \mathcal{B}(A))$ , such that  $d\bar{\pi}_\theta = \pi_\theta da$ .

$\Rightarrow$  (we consider  $X$  and  $A$  as subsets of Euclidean spaces)

- Consider discounted setting. (Ignoring a.s.)

$$V_{\pi_\theta} = E_{\pi_\theta} \left[ \sum_{t \geq 1} \gamma^{t-1} r_t \mid \sigma(x_1) \right]$$

$$Q_{\pi_\theta} = E_{\pi_\theta} \left[ \sum_{t \geq 1} \gamma^{t-1} r_t \mid \sigma(x_1, A_1) \right]$$

For a given  $x$

$$\Rightarrow V(x) = E_{\pi_\theta} \left[ Q(x, A_1) \mid \sigma(x_1) \right](x)$$

$$= \int_A \pi_\theta(a; x) Q_{\pi_\theta}(a; x) da$$

Therefore,

$$\begin{aligned} \nabla_{\theta} V_{\pi_{\theta}}(n) &= \nabla_{\theta} \int_A \pi_{\theta}(a; n) Q_{\pi_{\theta}}(a; n) da \\ &= \int_A \nabla_{\theta} \pi_{\theta}(a; n) Q_{\pi_{\theta}}(a; n) da \\ &\quad + \int_A \pi_{\theta}(a; n) \nabla_{\theta} Q_{\pi_{\theta}}(a; n) da \\ &\Rightarrow \bar{r}(n, a) + \gamma \int P(n'; n, a) V_{\pi_{\theta}}(n') dn' \end{aligned}$$

(we used the fact that  $d\bar{P}(x'; n, a) = P(x'; n, a) dn'$ ,  $n, n' \in \mathcal{X}$ ,  $a \in \mathcal{A}$ )

$$\begin{aligned} \nabla_{\theta} V_{\pi_{\theta}}(n) &= \int_A \nabla_{\theta} \pi_{\theta}(a; n) Q_{\pi_{\theta}}(a; n) da \\ &\quad + \gamma \int_A \pi_{\theta}(a; n) \left( \int_{\mathcal{X}} P(n'; n, a) \nabla_{\theta} V_{\pi_{\theta}}(n') dn' \right) da \end{aligned}$$

Let's define occupancy kernel; for  $S \in \mathcal{B}(\mathcal{X})$

$$\bar{\mu}_{\pi_{\theta}}^{\gamma}(S; n) = \lim_{T \rightarrow \infty} \frac{1}{1-\gamma} E_{\pi_{\theta}} \left( \sum_{t=1}^T \gamma^{t-1} I(X \in S) \mid \sigma(X_0) \right)(n)$$

$$\Rightarrow \mu_{\pi_{\theta}}^{\gamma} \text{ such that } d\bar{\mu}_{\pi_{\theta}}^{\gamma}(\cdot; n) = \mu_{\pi_{\theta}}^{\gamma}(n'; n) dn'$$

$$\text{Side note } V_{\pi_0}(n) = (1-\gamma) \int \int_{\mathcal{X} \times \mathcal{A}} \gamma \underbrace{\pi(a; n)}_{\pi_0(a; n)} \underbrace{r(n, a)}_{r(n, a)} \underbrace{\mu_{\pi_0}^{\gamma}(n', n)}_{\mu_{\pi_0}^{\gamma}(n', n)} da dn'$$

$$\nabla_{\theta} V_{\pi_0}(n) = \int \int_{\mathcal{X} \times \mathcal{A}} \gamma \underbrace{\pi_{\theta}(a; n) Q_{\pi_0}(a; n') \mu_{\pi_0}(n', n)}_{\pi_{\theta}(a; n) Q_{\pi_0}(a; n') \mu_{\pi_0}(n', n)} dn'$$

$$\Rightarrow \rho(\pi_0) = \mathbb{E}_{\pi_0} [\nabla_{\theta} V_{\pi_0}]$$

Undiscounted setting.

$$\rho(\pi_0) = \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\pi_0} \left[ \sum_{t=1}^T r_t \right]$$

$$\text{Let } \bar{\mu}_{\pi_0}(S) = \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E}_{\pi_0} \left[ \sum_{t=1}^T I(X \in S) \right]$$

$\hookrightarrow$  occupancy kernel.

$$S \in \mathcal{B}(\mathcal{X})$$

$$(\text{remember } \bar{P} = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T p^{t-1})$$

(we assume all the limits exist and is stationary)

$$\rho(\pi_0) = \int \int_{\mathcal{X} \times \mathcal{A}} \pi(a; n) \bar{r}(n, a) da \mu_{\pi_0}(n) dn$$

$$\text{where } d\bar{\mu}_{\pi_0} = \int_{\pi_0} dn$$

what is  $\nabla_{\theta} \rho(\pi_{\theta})$

remember

$$V_{\pi_{\theta}} = \lim_{T \rightarrow \infty} \frac{1}{T} E_{\pi_{\theta}} \left[ \sum_{t=1}^T r_t - \rho(\pi_{\theta}) \mid \sigma(x_1) \right]$$

$$\text{and } Q = \lim_{T \rightarrow \infty} \frac{1}{T} E_{\pi_{\theta}} \left[ \sum_{t=1}^T r_t - \rho(\pi_{\theta}) \mid \sigma(x_1, A_1) \right]$$

$$\Rightarrow V_{\pi_{\theta}}(x) = E_{\pi_{\theta}} [Q(x_1, A_1) \mid \sigma(x_1)](x)$$

$$\Rightarrow \int_{\mathcal{A}} Q_{\pi_{\theta}}(x, a) \pi_{\theta}(a; x) da$$

$$\Rightarrow \nabla_{\theta} V_{\pi_{\theta}}(x) = \int_{\mathcal{A}} \nabla_{\theta} \pi_{\theta}(a; x) Q_{\pi_{\theta}}(x, a) da$$

$$+ \int_{\mathcal{A}} \pi_{\theta}(a; x) \nabla_{\theta} Q_{\pi_{\theta}}(x, a) da$$

$$\rightarrow \underbrace{\bar{r}(x, a) - \rho(\pi_{\theta})}_{\text{}} + \underbrace{\int_{\mathcal{X}} \pi(x'; x, a) V_{\pi_{\theta}}(x') dx'}_{\text{}} \frac{d\pi_{\theta}}{d\theta}$$

$$\nabla_{\theta} V_{\pi_{\theta}}(n) = \int_{\mathcal{A}} \nabla_{\theta} \pi_{\theta}(a; n) Q_{\pi_{\theta}}(n, a) da$$

$$+ \int_{\mathcal{A}} \pi_{\theta}(a; n) \left( -\nabla_{\theta} \rho(\pi_{\theta}) + \int_{\mathcal{X}} p(n'; n, a) \nabla_{\theta} V_{\pi_{\theta}}(n') dn' \right)$$

$$\Rightarrow \nabla_{\theta} \rho(\pi_{\theta}) = \int_{\mathcal{A}} \nabla_{\theta} \pi_{\theta}(a; n) Q_{\pi_{\theta}}(n, a)$$

$$+ \int_{\mathcal{A}} \pi_{\theta}(a; n) \int_{\mathcal{X}} p(n'; n, a) V_{\pi_{\theta}}(n') dn' - \nabla_{\theta} V_{\pi_{\theta}}(n)$$

$$\Rightarrow \int \nabla_{\theta} \rho(\pi_{\theta}) \mu_{\pi_{\theta}} dn = \int_{\mathcal{X}} \mu_{\pi_{\theta}} \left( \int_{\mathcal{A}} \nabla_{\theta} \pi_{\theta}(a; n) Q_{\pi_{\theta}}(n, a) da \right) dn$$

$$+ \int_{\mathcal{X}} \mu_{\pi_{\theta}} \left( \int_{\mathcal{A}} \pi_{\theta}(a; n) \left( \int_{\mathcal{X}} p(n'; n, a) \nabla_{\theta} V_{\pi_{\theta}}(n') dn' \right) da \right) dn$$

$$- \int_{\mathcal{X}} \mu_{\pi_{\theta}} \nabla_{\theta} V_{\pi_{\theta}} dn$$

$$\Rightarrow \nabla_{\theta} \rho(\pi_{\theta}) = \int_{\mathcal{X}} \mu_{\pi_{\theta}} \left( \int_{\mathcal{A}} \nabla_{\theta} \pi_{\theta}(a; n) Q_{\pi_{\theta}}(n, a) da \right) dn$$

# Function, Measure, Kernel, Functional, Operator

$\rightarrow X \quad \sigma(X) \quad f: X \rightarrow Y \leftarrow$

$Y \quad \sigma(Y)$

$\mu; \sigma(X) \rightarrow \mathbb{R}$

new Markov kernel  $\Leftarrow \begin{cases} K(\cdot; \frac{1}{n}) \rightarrow \text{measure} \\ \text{on } \sigma(Y) \\ K(y; \cdot) \text{ measurable function} \end{cases}$

---

$\mathcal{H}$  is a Banach space e.g. space of function  
functional  $F: \mathcal{H} \rightarrow \mathbb{R}$

operator  $\mathcal{H}, G$

$T: \mathcal{H} \rightarrow G$

