# Week 7 - Reinforcement Learning

## GPI — Monday

grad ascent method

$\qquad \hookrightarrow \nabla$ for max

$\hookrightarrow \nabla$ of what?

$$\nabla V_\gamma^\ell (\pi) \overset{def}{=} \boxed{\frac{\partial V_\gamma^\ell (\pi)}{\partial \pi}} \rightsquigarrow \nabla / \text{partial derivative}$$

difficult to use / ~~identify~~

$\hookrightarrow \nabla$ gantinya / as a remedy $\underset{\text{policy parameteriz}}{\hookrightarrow \text{"explicit"}}$

$\nabla$ we introduce
a param vector
$\Theta \in \mathbb{R}^{dim(\theta)}$ , so $V(\pi(\theta)) \overset{def}{=} V(\theta)$

⟶ Best Practice

## Why explicit?

① can utilize high dimension Repr. for policy

② can encode levels of exploration

"share the same idea with $\varepsilon$-greedy"

③ Cari yang lain for "Exceeding esp".

---

## The Grad Ascend update rule:

Policy Gradient

$$\theta^{H+1} \leftarrow \theta^K + \alpha \nabla V_x^{\varrho}(\theta)\Big|_{\theta=\theta^K}$$

$\theta = \theta^K$
↳ evaluated using $\theta^K$

$K =$ policy iteration index

$\theta =$ initialized / choose randomly

$\alpha =$ choose random small number

$\nabla V_x(\theta) = ?$ ⟶ sutton's policy grad

set $\times \to \nabla$ gain

$$\nabla V_g(\theta) = \nabla \left\{ \sum_{s \in S} \sum_{a \in A} r(s,a) \, \underset{\pi(\theta)}{P^*(s)} \cdots \right.$$

$$\left. \cdots \pi(a|s;\theta) \cdots \text{ by def of gain} \right\}$$

$$= \sum_{s} \sum_{a} r(s,a) \left\{ \nabla P^*_\pi(s) \cdot \pi + P^*_\pi(s) \cdot \nabla \pi \right\}$$

$\nabla$ if we stop here, still its not sampling

friendly ~~tua~~ in RL , thus we use

"score function".

$$= \sum_{s} \sum_{a} r(s,a) \left\{ \left( P^*_\pi \cdot \nabla \log P^*_\pi \cdot \pi \right) + P^* \cdot \pi \cdot \nabla \log \pi \right\}$$

$$= \sum_{s} \sum_{a} P^*(s) \cdot \pi(a|s) \cdot r(s,a) \left\{ \nabla \log P^* + \nabla \log \pi \right\}$$

$\to$ sampling -friendly but in RL, $P^*$ is still

unknown to the agent.

Score function :

$$\nabla \log \pi (a|s;\theta) = \frac{\nabla \pi (a|s;\theta)}{\pi (a|s;\theta)}$$

$$= \sum_{s} \sum_{a} p^*(s) \, \pi(a|s) \, q_b^\pi (s,a) \, \nabla \log \pi(a|s;\theta)$$

$$= \mathbb{E} \left[ q_b^\pi (S|A) \, \nabla \log \pi(A|s;\theta) \right]$$

$S \sim p^*$

$A \sim \pi(\cdot|s)$

↳ exceeding exp
find proof!

$p^*$ is unknown but appears under $\mathbb{E}$, hence in practice we simply sample forom $P_\pi^* (\theta)$

↳ $p^*$ exist, at $t \geqslant t_{mix}$

run unbiased approx of $\nabla V(\gamma)$

but, at $t < t_{mix}$, $S \sim p^t \neq p^{\bigstar}$

hence biased approx. of $\nabla v(x)$

① . BCV : Reduce var w/o introducing bias error

$$\nabla V_g(\theta) = \mathbb{E}\left[\left\{q_b^\pi(S,A) - \overset{\text{baseline}}{\underbrace{V_b(S)}}\right\} \nabla \log \ldots\right.$$

$$\left.\ldots \pi(A|S;\theta)\right]$$

$$= \underbrace{\partial_b^\pi(S,A)}_{\text{advantage function}}$$

$$= \mathbb{E}\left[ \delta_v(S, A, S') \nabla \log \pi(A|S; \theta) \right]$$

# Week 7 — Reinforcement Learning

Wednesday.

Policy param for Discrete Action (focus: −
− finite MDP)

using: Gibbs / Boltzman / ~~softmax~~ categorical dist

softmax parameterization

$$\pi(a|s \; ; \theta) = \frac{\exp\left(\theta^T \; \phi(s,a)\right)}{\sum \exp\left(\theta^T \; \phi(s,a)\right)} \quad \forall (s,a) \in (S \times A)$$

# Discounted Reward Policy Gradient

Recall:
$$V_\gamma(\pi) = P_\pi^\gamma \; r_\pi \qquad V_g^c(\pi) = P_\pi^{\mathcal{D}} \; r_\pi$$

(← our symbol for Discounted)

Where
$$P_\pi^\gamma = \lim_{t_{mix} \to \infty} \sum_{t=0}^{t_{mix}-1} (\gamma \cdot P_\pi)^t \qquad P_\pi^{\mathcal{D}} \doteq (\text{see Prev})$$

$\llcorner$ is not a proper distribution

because $\sum_{s=0} P_\pi^\gamma (s|s_0) \neq 1$

$$\nabla V_{g\gamma}(\pi, s_0) = \sum_s P_\pi^\gamma (s|s_0) \sum_{a \in A} q_\gamma^\pi (s,a) \cdots$$

$$\cdots \nabla \log \pi(A_a | s; \theta)$$

not sampling friendly, not a distribution

times with ~~~~ $(1-\gamma)$

$$= \frac{1}{(1-\gamma)} \sum (1-\gamma) P_{\pi}^{\gamma} (\varsigma | \varsigma_0) \sum_{a} \cdots$$

↪ becomes a geometric Distribution

↪ to make it back to equality.

$$\pi_g^* = \underset{\pi \in \Pi}{\arg\max} V_g (\pi, \varsigma_0)$$

$$\pi_\gamma^* = \underset{\pi \in \Pi \; \varsigma_0 \sim \textcircled{P}}{\arg\max} \mathbb{E} \left[ V_\gamma (\pi, \varsigma_0) \right] \rightarrow \text{initial state Distribution}$$

$\pi_\gamma^\#$ depends on $\overset{\circ}{p}$

Call non-uniform optimal