

RL W-7 (17/3/2025).

Before Jumping Into Policy Grad - -

GP<sub>T</sub> → poleval → Given policy  $\pi$ ,  
what is  $V(\pi)$ ?

↓  
(Generalized)  
Inexact

polimprov → modify  $\pi^k$  such that  
 $V(\pi^{k+1}) \geq V(\pi^k)$   
"better"

$e_{PB}$   
 $e_{MS}$ : semi-grad → in RL → policy  
gradient.

P.g. via grad ascent method.  
↓  
of what?  
value. w.r.t its param  
 $DV_x(\pi) = \frac{\partial V_x(\pi)}{\partial \pi}$

To maximize i.e.,

$$\pi^* \leftarrow \arg\max_{\pi} V_x(\pi)$$

$\pi \in \Pi_S$

strictly "E"



FAKULTAS  
ILMU  
KOMPUTER

polgrad in direct policy parameterization  
hard to use → not beneficial

As a remedy for direct policy parameterization;

→ Use explicit policy parameterization!

How?

↓ where we introduce

param vector  $\theta \in \Theta = \mathbb{R}^{\dim(\theta)}$

So:  $V(\pi(\theta)) \stackrel{\text{def}}{=} V(\theta)$ .

Example:

Let  $\pi \in \Pi_{SR}$

proven  
empirical  
success

Set  $\pi \sim \text{Gaussian}(\mu=0, \sigma=\theta)$ .

Why? Now it's becoming best practice & SOTA!

- ① utilize high-dimension representation for policy  
e.g., nets → learning  $\sigma$
- ② encode levels of exploitation  
③ ?  
↳ share the similar idea with  $\epsilon$ -greedy.



The Gradient Ascent Rule (for inexact policy improvement)

$$\theta^{k+1} = \theta^k + \alpha \cdot \nabla V_x(\theta) \quad | \quad \theta = \theta^k$$

positive step size

Assume:

$V_x(\theta)$  is differentiable

evaluated using  $\theta^k$

iteration index

Example:

$$f(x) = 2x^2 \rightarrow f'(x) = 4x \Big|_{x=1}$$

How?

$\theta$ : init, e.g., randomly

require smoothness

↓  
stochastic!

$$f(1) = 4,$$

$\alpha$ : init e.g.,  $10^{-3}$

deterministic  $\rightarrow$  Not differentiable

Big Question

Answer:

Sutton's Policy Gradient Theorem:

$$\nabla V_g = P^\star \left[ V_g = P \cdot r_R \right] \rightarrow \text{matrix form}$$

Set  $x \leftarrow \text{gain}$ .

$$\nabla V_g = \nabla \left\{ \sum_{S \in S} \sum_{a \in A} r(s, a) P^\star(s) \Pi(a|s; \theta) \right\} \quad \} (\text{by def of gain}).$$



FAKULTAS  
ILMU  
KOMPUTER

$$= \sum_s \sum_a r(s, a) \nabla \left\{ P_{\Pi(\theta)}^\star(s) \Pi(a|s; \theta) \right\}$$

Note that  $P_{\Pi(\theta)}^\star$  ~~is not~~ depends on  $\theta$ , as  $\Pi$  depends on  $\theta$

Continuing on  $\nabla V_g \dots$   
 Apply derivative's product rule

$$= \sum_s \sum_a r(s, a) \left\{ \nabla_{p_n^*} \cdot \pi + p_n^* \cdot \nabla \pi \right\}.$$

↓

We don't know  $p_n^*$  → Not sampling friendly!

Thus, we need score function!  
 ↓.

$$\nabla \log \pi(a|s; \theta) = \frac{\nabla \pi(a|s; \theta)}{\pi(a|s; \theta)} \quad \begin{cases} \text{diff of} \\ \text{log.} \end{cases}$$

$$= \sum_s \sum_a r(s, a) \left\{ p_n^* \nabla \log p_n^* \pi + p_n^* \cdot \pi \cdot \nabla \log \pi \right\}$$

(by score function)

~~$$= \sum_s \sum_a p^*(s) \pi(a|s) \cdot r(s, a) \left\{ \nabla \log p^* + \nabla \log \pi \right\}$$~~

Sampling friendly, but in RL  
 $p^*$  is unknown to the agent!

Fun fact: Sutton spent days to solve this.



Cont... . There's  $p^*$  still, but it's becoming sample distribution  $s \sim p^*(\cdot)$ .

$$= \sum_s \sum_a p^*(s) \cdot \pi(a|s) \cdot q_b^n(s, a) \cdot \nabla \log \pi(a|s; \theta)$$

$\nabla$

Encode  $\nabla \log p^*$  and  $r(s, a)$ .

$$= E [q_b^n(s, a) \cdot \nabla \log \pi(a|s; \theta)] \quad \text{proof?}$$

$s \sim p^*(\cdot)$

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

unknown but appears under  $E$

RV.

a.k.a. mean  $\bar{m} \approx$  sample mean,  $\bar{m}$ , which is unbiased, i.e.,  $E[\bar{m}] = m$ .

hence stochastic gradient ascent

sample mean

unbiased

run unbiased approx of  $\nabla V_{\pi}(\theta)$   
hence SGA

hence in practice, we simply sample from  $(p^*, \pi(\cdot|s))$  to exist at  $t \geq t_{mix}$



As a programmer, likely  $t_{\text{mix}}$  can be computed in small envs, but not in large scale, thus it's becoming hyperparams!

At  $t < t_{\text{mix}}$ ,  $s \approx p^+ \neq p^*$ , hence it's biased approx of  $\nabla V(\pi)$ .  
↳ semi-SGA

Code:

Gymnium  $\rightarrow$  env library  
stable-baseline  $\rightarrow$  algo library -

We can also  
Create custom  
Env's!



Cont. . .

Issue: Sample mean of  $E[q \nabla \log \pi]$  is unbiased but "generally" has high variance as:

involves  $p^*$  and  $\Omega \rightarrow 2$  RVs!

as a remedy . . . ~~fat baseline~~, e.g.,

- ① Baseline control variate (BCV) can complement each other!
- ② Actor-critic

→ ① BCV: goal: reduce variance without introducing ~~bias-error~~ different from b in  $V_b$ .

How: use a baseline

independent of action; e.g.,  $V_b(s)$

Hence,

advantage (state-action) function ~~function~~

$$\nabla V_g(\theta) = E \left[ \left\{ q_b^n(s, a) - V_b(s) \right\} \nabla \log \pi(a|s; \theta) \right]$$

~~$= E[\delta_n(s, a)]$~~  baseline

Fun fact :

⇒ theorem:

$$\hat{J}(s, a) = \mathbb{E} \left[ r(s, a) - V_g + V_b(s') - V_b(s) \mid s, a \right]$$

$s' \sim p(\cdot | s, a)$       TD :  $\delta_v$

Says TD is ~~an~~ unbiased estimate for  $\hat{J}$ !

Thus :

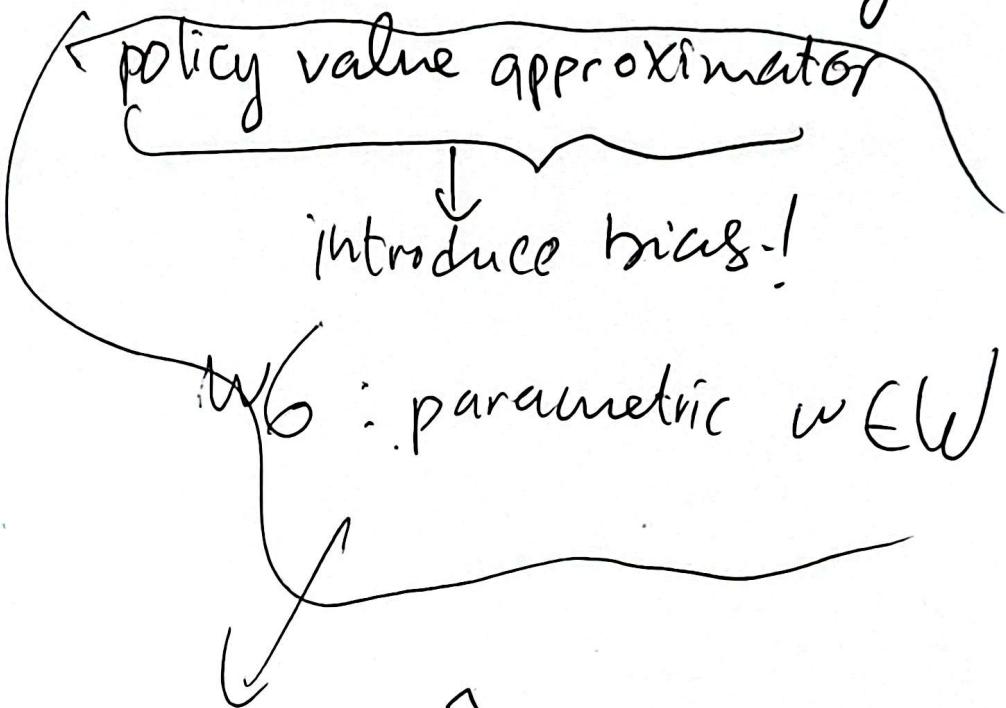
$$= \mathbb{E} \left[ \delta_v(s, a, s') \nabla \log \pi(a|s; \theta) \right]$$

Using sample of  $(s, a, s')$ :

$$\approx \delta_v(s, a, s') \nabla \log \pi(a|s; \theta)$$



② Actor Critic: reduce variance using



a.k.a Critic!  $\xrightarrow{\text{forms}} V_b(s; w)$

actor refers to:  ~~$f(a|s; \theta)$~~ ,  
 $\pi(\theta)$

Here we  
apply  
both:

① baseline.

which used in



$\delta_t$  for approx  $d$   
for approx  $DV_g(\theta)$



FAKULTAS  
ILMU  
KOMPUTER

② Actor critic

# RL Week -7 ~2

Policy Param for Discrete Action.

(Focus : finite MDP).

↳  $S, A$  finite

e.g., using Gibbs/Boltzmann/Categorical Distribution

Softmax.

$$\pi(a|s; \theta) = \frac{\exp(\theta^T \phi(s, a))}{\sum_{a' \in A} \exp(\theta^T \phi(s, a'))}; \quad \forall (s, a) \in S \times A$$

↓  
conditional probability  
over action.

Recall:

$$\pi: S \rightarrow \Delta^{(A)}$$

or ~~or~~

$$\pi: S \times A \rightarrow [0, 1].$$

e.g.,  $s^0 \rightarrow \boxed{0.3 \mid 0.6 \mid 0.1}$



# Discounted Reward - Policy Gradient -

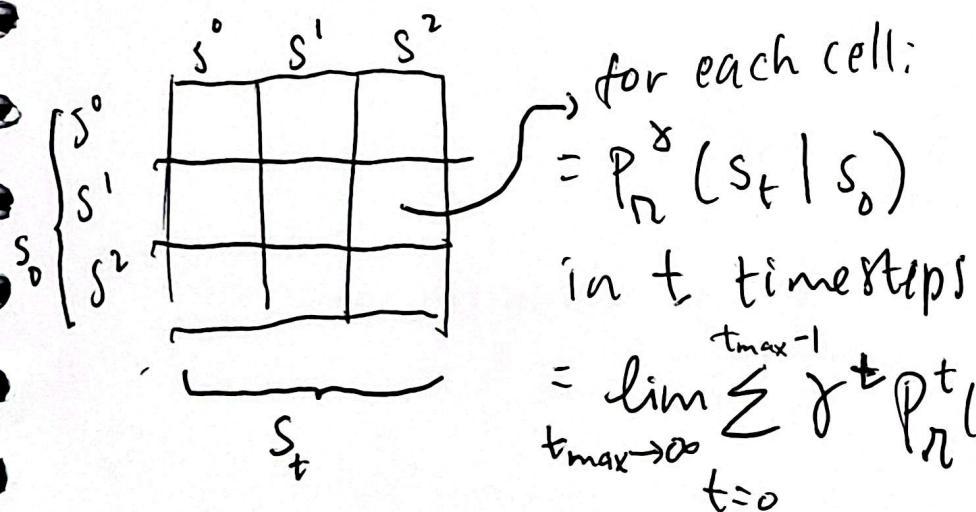
Recall that :

$$V_g(\pi) = P_n^* r_n \quad \text{where } P_n^* = \lim_{t \rightarrow \infty} P_n^t$$

$$V_\gamma(\pi) = P_n^\gamma \cdot r_n$$

where

$$P_n^\gamma = \lim_{t_{\max} \rightarrow \infty} \sum_{t=0}^{t_{\max}-1} (\gamma \cdot P_n)^t$$



is  $P_n^\gamma$  a distribution?

If we sum the elements in a row, it'll be 1, e.g.

$$\sum_{s \in S} \{P_n^\gamma(s | s_0)\} = 1.$$

$$\sum_{s \in S} \left\{ \lim_{t_{\max} \rightarrow \infty} \sum_{t=0}^{t_{\max}-1} \gamma^t \cdot P_n^t(s | s_0) \right\}$$



FAKULTAS  
ILMU  
KOMPUTER

$P_\pi^\gamma$  is improper distribution, as  $\frac{1}{1-\gamma} \neq 1$

Then, how's next?

Recall that:-

$V_\gamma$  depends on initial state  $\rightarrow$

introduce complexity  
c.f.  $\nabla V_\theta(\pi)$

$$\nabla V_\gamma(\pi, s_0) = \sum_s P_\pi^\gamma(s|s_0) \sum_{a \in A} q_\gamma^n(s, a) \cdot \nabla \pi(a|s; \theta)$$

not sampling friendly

as  $P_\pi^\gamma$  is not a ~~distribution~~ distribution.

Thus, we introduce  $\frac{1-\gamma}{1-\gamma}$ :

$$\nabla V_\gamma(\pi, s_0) = \frac{1}{1-\gamma} \sum_s \underbrace{(1-\gamma) P_\pi^\gamma(s|s_0)}_{\text{proper distribution!}} \sum_{a \in A} q_\gamma^n(s, a) \cdot \nabla \pi(a|s; \theta)$$

proper distribution!

i.e. Geometric Distribution  
 $(p = 1-\gamma)$ .

Q: How to sample from that distribution?



INSTITUT  
ILMU  
KOMPUTER

We can also scale the gradient;

just like ~~what~~ what  $\alpha$  (learning rate) do:

$$(1-\gamma) \cdot \nabla V_\gamma(\pi, s_0) = \sum_s (1-\gamma) P_\pi^\gamma(s|s_0) \sum_{a \in A} \dots$$

Difference: ① State Distribution:

$$S \sim p_n^*$$

avg rew

dist. rew.

$$S \sim \text{Geo}(p=1-\beta)$$

number of trials  
until success, i.e. ~~terminates~~  
terminates after  $t+1$ .

② Dependency on  $S_0$

avg rew

$$\pi^* = \underset{\pi \in \Pi}{\operatorname{argmax}} . V_g(\pi)$$

uniform optimal  
w.r.t.  $\hat{p}$ .

discounted rew

$$\pi_g^* = \underset{\pi \in \Pi}{\operatorname{argmax}} \mathbb{E} [V_g(\pi, S_0)]$$

init  
state  
dist

Interpretation:

$\pi_g^*$  depends on  $\hat{p}$  ~~→ non~~

non-uniform optimal

w.r.t.  $\hat{p}$ , e.g.;

$$\hat{P}_{(1)} = \begin{array}{|c|c|c|} \hline s^0 & s^1 & s^2 \\ \hline 1 & 0 & 0 \\ \hline \end{array}$$

$$\hat{P}_{(2)} = \begin{array}{|c|c|c|} \hline & & \\ \hline 0 & 1 & 0 \\ \hline \end{array}$$



How To Sample From  $S \sim \text{Geo}(p=1-\gamma)$   
 See Paper . . . —

$$\nabla V_\gamma(\pi, \xi_0) = \frac{1}{1-\gamma} \sum_{s \in S} ((1-\gamma) \cdot p_n^\gamma \sum_{a \in A} \dots \\ = \frac{1}{1-\gamma} \mathbb{E} [ \dots ].$$

$$S \sim (1-\gamma) \cdot p_n^\gamma \cdot$$

$$A \sim \pi.$$

How to sample states from  $(1-\gamma) \cdot p_n^\gamma = \tilde{p}_n^\gamma$ ?  
 $\text{Geo}(p=1-\gamma)$

① Sample a length of xprmt-episode from  $\text{Geo}(1-\gamma)$   
num of steps = length(xprmt-episode).

$$L \sim \text{Geo}(1-\gamma).$$

Say  $L = 20$  steps

② Run an xprmt-episode  $t=0 \rightarrow t=20-1$ .

③ Thm  $S_{t=19} \sim \tilde{p}_n^\gamma \rightarrow$  we throw away all state samples from  $t=0$  to  $t=18$

In practice, people use all state from  $\tilde{p}_n^\gamma \rightarrow$  Not sample efficient!  
 → biased → semi-SGA!

