

Week-10 - RL

1 Reinforce : MC Policy Gradient

$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a q_{\pi}(s, a) \nabla J(\alpha | s, \theta)$$

$$\theta_{t+1} = \theta_t + \alpha G_t \frac{\nabla J(A_t | s_t, \theta_t)}{J(A_t | s_t, \theta_t)}$$

- it uses the policy of comp. return from time t , which includes all the future rewards up until the end of the episode.
- the algo. computes an unbiased estimate of the gradient. Can be very slow due to high variance (occurs due to recurrence time).
- With Baseline

$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a (q_{\pi}(s, a) - b(s)) \nabla J(\alpha | s, \theta)$$

Update rule ,

$$\theta_{t+1} = \theta_t + \alpha (G_t - b(s_t)) \frac{\nabla J(A_t | s_t, \theta_t)}{J(A_t | s_t, \theta_t)}$$

12 Actor Critic Methods - I

- ACM \rightarrow learn both a policy & a state value funcⁿ simultaneously.
The policy is π , referred to as the actor that suggests action given a state.
- The estimated value funcⁿ is called the critic. It evaluates the ac^{"s} taken by the actor based on the given policy.

$$\theta_{t+1} = \theta_t + \alpha \left(G_{t:t+1} - \hat{v}(s_t, w) \right)$$
$$\frac{\Delta \pi(A_t | s_t, \theta_t)}{\pi(A_t | s_t, \theta_t)}$$

Introducing TD error,

$$\theta_{t+1} = \theta_t + \alpha S_t \frac{\nabla \pi(A_t | s_t, \theta_t)}{\pi(A_t | s_t, \theta_t)}$$

- 1 step actor critic method, is like an episodic task for estimating the $\pi_\theta \approx \pi$.
- comparing the current method to baseline

1. The G_t term although unbiased causes high variance.
2. since, $E_{\pi}[G_t | S_t, A_t] = q_{\pi}(S_t, A_t)$
if somehow, we had an estimate of $q_{\pi}(S_t, A_t)$ with less variance then we can use it instead of G_t .
3. In 1 step AC, we use \hat{v} for both estimation of $q_{\pi}(S_t, A_t)$ & the base
4. The bootstrapping method in the update introduces bias but reduces variance.
5. reduced variance helps in increased learning.

L3 ACM - 2

- actor - computes the policy π_θ and updates θ .
- critic - computes the estimate $\hat{v}(s, w)$ of the state value funcⁿ. Also, updates the parameter w .
- basic algo. for actor critic -
 1. take acⁿ $a \sim \pi_\theta(a|s)$ & receives (s, a, s', r)
 2. update w using data $(s, r + \gamma \hat{v}(s', w))$

3. calculate,

$$\hat{\delta}(s, a) = r + \gamma \hat{v}(s', w) - \hat{v}(s, w)$$

4. $\theta \leftarrow \theta + \alpha \cdot \hat{\delta}(s, a) \cdot \nabla_{\theta} \log \pi_{\theta}(a|s)$

5. repeat

- how to update the critic?

→ we will have multiple data points
of the form (s, y)

→ min. the squared loss

$$L(w) = \frac{1}{N} \sum_i \| \hat{v}(s_i, w) - y_i \|^2$$

↳ is the batch size

- Advantage funcⁿ → the diff. b/w q
value & value funcⁿ.

$$A_{\pi}(s, a) = q_{\pi}(s, a) - v_{\pi}(s)$$

- A3C → asynchronous advantage actor
critic. Training happens parallelly

- A2C → synchronous adv. actor critic.

↳ presence of a coordinator. The
updates of the global parameters are
executed only after the threads have
finished their computaⁿ.

L4 DPG & DB PG \longrightarrow deep DPG

Deterministic policy gradient

- policy gradient over continuous spaces

$$\nabla J(\theta) = E_{\pi} \left[\sum_a q_{\pi}(s_t, a) \nabla \pi(a | s_t, \theta) \right]$$

\hookrightarrow for discrete ac^{"s}

$$\nabla_{\theta} J(\theta) = E_{s \sim \mu^{\pi}, a \sim \pi_{\theta}} \left[\nabla_{\theta} \log \pi_{\theta}(a | s) q_{\pi_{\theta}}(s, a) \right]$$

\hookrightarrow for continuous ac^{"s}

- it is often hard to implement the differentiable continuous controllers in many problems. Expecta["] over both states & ac^{"s}.

- deterministic policies

$$\pi_{\theta} : s \rightarrow A \quad (\text{det.})$$

$$J(\theta) = \int_S \mu^{\pi}(s) q_{\pi_{\theta}}^{\theta}(s, \pi_{\theta}(s)) ds$$

\hookrightarrow the continuous ac["] spaces allow us to think of θ in ac["] w.r.t the policy parameters. Now, we will have expecta["] only over states.

- DPG - it is a limiting case of stochastic policy gradient for a very wide class of stochastic policies.
- If the env. is stochastic, then not a problem. Otherwise, we will use an off-policy actor critic, where the behavioural policy differs from the "target" policy.
- DDPG = DPG + DQN (off-policy)
$$\pi'(s) = \pi_\theta(s) + N \text{ (noise)}$$
- use replay buffer
- maintains separate target network parameters θ' , w' & uses soft updates.
- use batch normalization method to normalize input state features & minimize covariate shift.