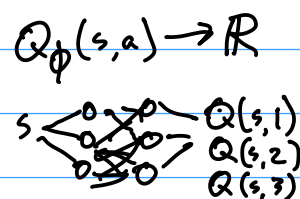


HW 4 - 1 alg from notebook
1 alg implemented

Last Time

DQN \leftarrow Experience Replay
Freezing target Q network
NN Architecture



DPG \sim Natural Gradient

Rainbow

Today

Actor-Critic

Exploration

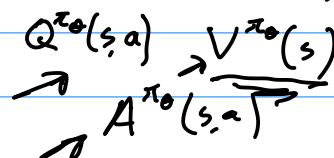
Recap of RL

Actor-Critic

$\pi_\theta \rightarrow \theta$ $Q_\phi \rightarrow \phi$ or A_ϕ or V_ϕ

$$\nabla_\theta U(\theta) = E \left[\sum_{k=1}^d \nabla_\theta \log \pi_\theta(a^{(k)} | s^{(k)}) \gamma^{k-1} \left(\overbrace{r_{to go}^{(k)} - r_{base}(s^{(k)})} \right) \right]$$

$$A^{\pi_\theta}(s,a) = E[r + \gamma V^{\pi_\theta}(s') - V^{\pi_\theta}(s)]$$



$$\ell(\phi) = \frac{1}{2} E_s \left[(V_\phi(s) - V^{\pi_\theta}(s))^2 \right]$$

GAE: Generalized Advantage Estimation

$$A^{\pi_\theta}(s,a) \leftarrow E[r + \gamma V_\phi(s') - V_\phi(s)] \quad \leftarrow \text{Bias}$$

$$E[r^k + \gamma r^{(k+1)} + \gamma^2 r^{(k+2)} + \dots + \gamma^{d-1} r^{(d)} - r_{base}(s)]$$

$$E[r^{(k)} + \gamma r^{(k+1)} + \dots + \gamma^m V^{\pi_\theta}(s^m) - V^{\pi_\theta}(s)] \quad \leftarrow \text{Variance}$$

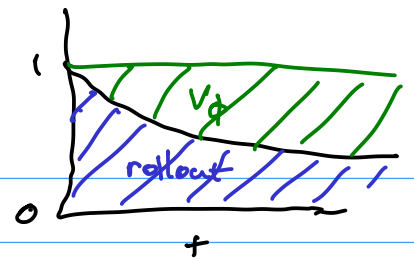
$$\delta_t = r_t + \gamma V_\phi(s_{t+1}) - V_\phi(s_t) = \hat{A}_t$$

$$\hat{A}^{(k)}(s,a) = E \left[\sum_{l=1}^k \gamma^{(l-1)} \delta_l \right]$$

$$\hat{A}^{GAE}(s,a) = E \left[\sum_{l=1}^k (\gamma \lambda)^{(l-1)} \delta \right]$$

$$\lambda \in [0,1]$$

$\lambda = 1$: rollout
 $\lambda = 0$: value



Acronyms to know

A3C Asynchronous Advantage Actor Critic

SAC $\pi(a|s) \propto e^{Q(s,a)}$

$$\pi_{\text{MaxEnt}}^* = \operatorname{argmax} E \left[\sum_t r_t + H(\pi(a|s_t)) \right]$$

↑ Entropy

Continuous A

$$\pi(a|s)$$

DPG Deterministic Policy Gradient

$$a = \pi_{\theta}(s) \quad Q_{\phi}(s,a)$$

training $Q \rightarrow$ $l(\phi) = \frac{1}{2} E \left[(r + \gamma Q_{\phi}(s', \pi_{\theta}(s')) - Q_{\phi}(s,a))^2 \right]$

train π_{θ}

training $\theta \rightarrow U(\theta) = E_s [Q_{\phi}(s, \pi_{\theta}(s))]$

$$\nabla U(\theta) = E [\nabla_{\theta} Q_{\phi}(s, \pi_{\theta}(s))]$$

$$= E \left[\nabla_{\theta} \pi_{\theta}(s) \nabla_a Q_{\phi}(s,a) \Big|_{a=\pi_{\theta}(s)} \right]$$

↑
new

Exploration? add Gaussian Exploration

Alpha-Zero

Actor-Critic with MCTS

Learn $\pi_{\theta}(a|s)$

$V_{\phi}(s)$

Guide Tree Search

replace rollouts

$$a = \operatorname{argmax} Q(s,a) + c \pi_{\theta}(a|s) \frac{\sqrt{N(s)}}{1+N(s,a)}$$

train θ

$$\pi_{\text{MCTS}}(a|s) \propto N(s,a)^{\eta}$$

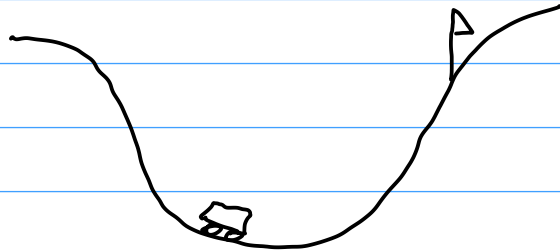
Cross entropy loss $l(\theta) = -E_s \left[\sum_a \pi_{\text{mcts}}(a|s) \log \pi_{\theta}(a|s) \right]$

train ϕ

$$V_{MCTS}(s) = \max_a Q(s,a)$$

$$\ell(\phi) = \frac{1}{2} \mathbb{E}_s \left[(V_\phi(s) - V_{MCTS}(s))^2 \right]$$

Advanced Exploration



Montezuma's Revenge

Breakout Rooms

How would you make an RL algorithm to handle MR

- Add incremental Rewards "Reward Shaping"
- Start state closer to goal
- Prune MCTS early
- Depth in Pyramid - Remembering

$$R^+(s,a) = \overset{\text{extrinsic}}{R(s,a)} + \overset{\text{intrinsic}}{B(s,a)}$$

Continuous States

fit $p_\theta(s)$ "pseudo-count"

What Bonus

$$\text{UCB} \quad B(s,a) = \frac{\sqrt{\log N(s)}}{\sqrt{N(s,a)}}$$

$$B(s) = \frac{\sqrt{T}}{\sqrt{N(s)}}$$

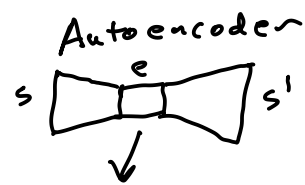
$$B(s) = \frac{1}{N(s)}$$

"Unifying Count-Based Exp."

Bellamare et al.

- "hash"
- # Exploration

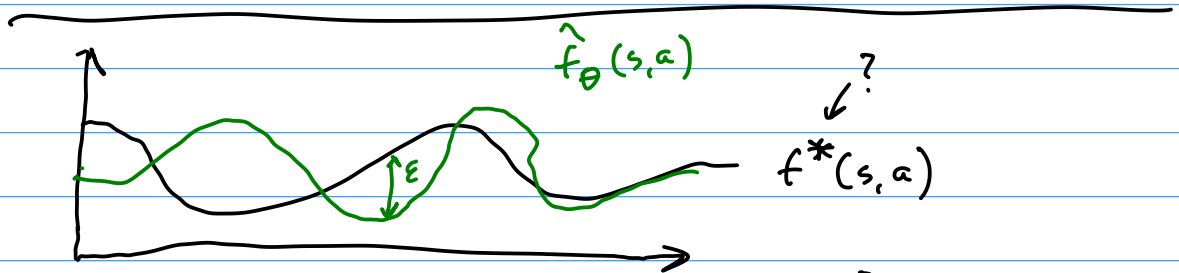
compress s into k -bit code $\phi(s)$ then count $N(\phi(s))$



- Use Classifier "EX2"

$$p_{\theta}(s) = \frac{1 - D_s(s)}{D_s(s)}$$

$D_s(s)$ prop that state returns true



use $\mathcal{E}(s, a) = \|\hat{f}(s, a) - f^*(s, a)\|^2$ as bonus

$\rightarrow f^*(s, a) = s'$ "Curiosity"

$\rightarrow f^*(s, a) = f_{\phi}(s, a)$ where ϕ is a random NN

Good \rightarrow "Random Network Distillation"

Thompson Sampling

Maintain distribution over $Q \leftarrow$ hard

Sample \hat{Q}_n

$a = \text{argmax } \hat{Q}(s, a)$

- Bootstrapping
- Multiple Q -networks

- Dropout

Information Gain

VIME

Go - Explore

First Return, Explore

Recap of RL

Exploration/Exploitation - Bandits, Advanced Exp.

Credit Assignment - Value Function, Eligibility/GAE

Generalization - Neural Network

Model Based

MLMBRL

BAMDP



Model Free

Learn Q

on policy

SARSA

off policy

Q-learning

Learning π

PG

Actor Critic