

$\hat{Q}_{i,t}$: estimate of expected reward - "reward to go", can be optimized,

↳ ideally: $\sum_{t'=t}^T E_{\pi_0} [r(s_{t'}, a_{t'}) | s_t, a_t]$ (true expected reward, all possible rewards)

total reward from taking a_t in s_t

$$\nabla_{\theta} J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_{i,t} | s_{i,t}) (\hat{Q}_{i,t}^{+}(a_{i,t}) - V^{\pi}(s_{i,t}))$$

better estimate lower variance. → "advantage" → $A^{\pi}(s_{i,t}, a_{i,t})$ baseline "average reward" $V(s_{i,t}) = E_{a_t \sim \pi_{\theta}}(a_t | s_t)$ $[Q(s_t, a_t)]$

→ essence: separately fit $Q^{\pi}/V^{\pi}/A^{\pi}$. (value-function network also actor-critic).

Bellman equation →

$$Q^{\pi}(s_t, a_t) = \underbrace{r(s_t, a_t)}_{\text{current reward}} + E_{s_{t+1} \sim p(s_{t+1} | s_t, a_t)} [\underbrace{V^{\pi}(s_{t+1})}_{\text{expectation of the future}}] \left(\sum_{t'=t}^T E_{\pi_0} [r(s_{t'}, a_{t'}) | s_t, a_t] \right)$$

$$\approx r(s_t, a_t) + V^{\pi}(s_{t+1})$$

$$A^{\pi}(s_t, a_t) \approx r(s_t, a_t) + V^{\pi}(s_{t+1}) - V^{\pi}(s_t) \leftarrow \text{just fit } V^{\pi}(s), \text{ neural network is separately trained } (s \rightarrow V) \text{ (policy evaluation)}$$

Policy Evaluation (fit V^{π})

generate samples (run the policy) → fit model → improve policy $Q \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

"Monte Carlo" method $V^{\pi}(s_t) \approx \sum_{t'=t}^T r(s_{t'}, a_{t'})$ a trajectory, $y_{i,t}$ "reward to end"

training data: $\{(s_{i,t}, \sum_{t'=t}^T r(s_{i,t'}, a_{i,t'}))\}$

- two network design

- shared network design

Two networks:

"Actor": policy π_{θ}

"Critic": value V^{π}

w/ Monte-carlo or bootstrap

gradient ascent

Example: TD-gammon (1992)

↳ reward: game outcome

↳ value function ($\hat{V}_{\theta}^{\pi}(s_t)$): expected outcome given board state

use regression + MSE. bootstrap $y_{i,t} = r(s_{i,t}, a_{i,t}) + V_{\theta}^{\pi}(s_{i,t+1})$

can improve: $y_{i,t} = r(s_{i,t}, a_{i,t}) + V_{\theta}^{\pi}(s_{i,t+1})$ previous fitted value function

batch actor-critic algorithm.

1. sample $\{s_i, a_i\}$ from $\pi_{\theta}(a|s)$

2. fit $\hat{V}_{\theta}^{\pi}(s_i)$ to sampled reward sums

3. evaluate $\hat{A}^{\pi}(s_i, a_i) = r(s_i, a_i) + \hat{V}_{\theta}^{\pi}(s_i) - \hat{V}_{\theta}^{\pi}(s_i)$

4. $\nabla_{\theta} J(\theta) \approx \sum_i \nabla_{\theta} \log \pi_{\theta}(a_i | s_i) \hat{A}^{\pi}(s_i, a_i)$

5. $\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta)$

discount factor γ : rewards sooner than later

$0 < \gamma < 1$, typically close to 1. can affect

$y_{i,t} \approx r(s_{i,t}, a_{i,t}) + \gamma \hat{V}_{\theta}^{\pi}(s_{i,t+1})$

MDD (ex. 18.8).

probably better than $\pi(a_t, s_t)$.

Q-learning: RL without policy gradient.

$\arg \max_{a_t} A^{\pi}(s_t, a_t)$: best action from s_t . → $\pi'(a_t | s_t) = \begin{cases} 1 & \text{if } a_t = \arg \max_{a_t} A^{\pi}(s_t, a_t) \\ 0 & \text{otherwise} \end{cases}$

fit A^{π} (or Q^{π}/V^{π}) (or estimate V^{π}).

generate samples

$\pi \leftarrow \pi'$

policy iteration algorithm.

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

$V^{\pi}(s) \leftarrow r(s, \pi(s)) + \gamma E_{s' \sim p(s' | s, \pi(s))} [V^{\pi}(s')]$

one caveat: curse of dimensionality.

→ use NN. 1. $y_i = \max_{a_i} (r(s_i, a_i) + \gamma E[V_{\theta}(s'_i)])$

2. $\theta \leftarrow \arg \min_{\theta} \frac{1}{2} \sum_i \|V_{\theta}(s'_i) - y_i\|^2$

"epsilon-greedy" exploration: $\pi(a_t | s_t) = \begin{cases} 1 - \epsilon & \text{if } a_t = \arg \max_{a_t} Q_{\theta}(s_t, a_t) \\ \epsilon / (|\mathcal{A}| - 1) & \text{otherwise} \end{cases}$

"Boltzmann exploration" $\pi(a_t | s_t) \propto \exp(Q_{\theta}(s_t, a_t))$ proportional.

can use replay buffer to sample i.i.d. (similar to batch)

1. take some action a_i and observe (s_i, a_i, s'_i, r_i) , add to \mathcal{B}

2. sample mini-batch $\{s_i, a_i, s'_i, r_i\}$ from \mathcal{B} uniformly

3. compute $y_i = r_i + \gamma \max_{a_i} Q_{\theta}(s'_i, a_i)$ using target network $Q_{\theta'}$

4. $\theta \leftarrow \theta - \alpha \sum_j \frac{\partial Q_{\theta}(s_j, a_j)}{\partial \theta} (Q_{\theta}(s_j, a_j) - y_j)$; 5. update θ' : copy θ every N steps

parameters: dataset size N , iterations K , gradient steps S

leave out $V^{\pi}(s)$ (unaffected)

larger buffer more stability.

gradient descent

"classic" QON algorithm

(use in HW#4)

$k=1$