Carnegie Mellon School of Computer Science

Deep Reinforcement Learning and Control

Maximum Entropy Reinforcement Learning

CMU 10-403

Katerina Fragkiadaki



RL objective

$$\pi^* = \arg\max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t} R(s_t, a_t) \right]$$

$$\pi^* = \arg\max_{\pi} \mathbb{E}_{(s_t, a_t) \sim \rho_{\pi}} \left[\sum_{t} R(s_t, a_t) \right]$$

MaxEntRL objective

Promoting stochastic policies

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t=1}^{T} \frac{R(s_t, a_t) + \alpha \operatorname{H}(\pi(\cdot \mid s_t))}{\operatorname{reward}} \right]$$

Why?

- Better exploration
- Learning alternative ways of accomplishing the task
- Better generalization, e.g., in the presence of obstacles a stochastic policy may still succeed.

Principle of Maximum Entropy

Policies that generate similar rewards, should be equally probable. We do not want to commit.

Why?

- Better exploration
- Learning alternative ways of accomplishing the task
- Better generalization, e.g., in the presence of obstacles a stochastic policy may still succeed.

We have seen this before.

Algorithm S3 Asynchronous advantage actor-critic - pseudocode for each actor-learner thread.

```
// Assume global shared parameter vectors \theta and \theta_v and global shared counter T = 0
// Assume thread-specific parameter vectors \theta' and \theta'_v
Initialize thread step counter t \leftarrow 1
repeat
     Reset gradients: d\theta \leftarrow 0 and d\theta_v \leftarrow 0.
     Synchronize thread-specific parameters \theta' = \theta and \theta'_v = \theta_v
     t_{start} = t
     Get state s_t
     repeat
          Perform a_t according to policy \pi(a_t|s_t; \theta')
          Receive reward r_t and new state s_{t+1}
          t \leftarrow t + 1
          T \leftarrow T + 1
     until terminal s_t or t - t_{start} = t_{max}
    R = \begin{cases} 0 & \text{for terminal } s_t \\ V(s_t, \theta'_v) & \text{for non-terminal } s_t // \text{Bootstrap from last state} \end{cases}
     for i \in \{t - 1, ..., t_{start}\} do
          R \leftarrow r_i + \gamma R
          Accumulate gradients wrt \theta': d\theta \leftarrow d\theta + \nabla_{\theta'} \log \pi(a_i | s_i; \theta') (R - V(s_i; \theta'_v) + \beta \nabla_{\theta'} H(\pi(s_t; \theta')))
          Accumulate gradients wrt \theta'_v: d\theta_v \leftarrow d\theta_v + \partial \left(R - V(s_i; \theta'_v)\right)^2 / \partial \theta'_v
     end for
     Perform asynchronous update of \theta using d\theta and of \theta_v using d\theta_v.
until T > T_{max}
```

"We also found that adding the entropy of the policy π to the objective function improved exploration by discouraging premature convergence to suboptimal deterministic policies. This technique was originally proposed by (Williams & Peng, 1991)"

MaxEntRL objective

Promoting stochastic policies

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\pi} \left[\sum_{t=1}^{T} \frac{R(s_t, a_t) + \alpha \operatorname{H}(\pi(\cdot \mid s_t))}{\operatorname{reward}} \right]$$

How can we maximize such an objective?

Recall:Back-up Diagrams



$$q_{\pi}(s,a) = r(s,a) + \gamma \sum_{s' \in \mathcal{S}} T(s' \mid s,a) \sum_{a' \in \mathcal{A}} \pi(a' \mid s') q_{\pi}(s',a')$$

Back-up Diagrams for MaxEnt Objective



Back-up Diagrams for MaxEnt Objective



$$q_{\pi}(s,a) = r(s,a) + \gamma \sum_{s' \in \mathcal{S}} T(s' \mid s,a) \sum_{a' \in \mathcal{A}} \pi(a' \mid s') \left(q_{\pi}(s',a') - \log(\pi(a' \mid s')) \right)$$

(Soft) policy evaluation

Bellman backup equation:

$$q_{\pi}(s,a) = r(s,a) + \gamma \sum_{s' \in \mathcal{S}} T(s' \mid s, a) \sum_{a' \in \mathcal{A}} \pi(a' \mid s') q_{\pi}(s', a')$$

Soft Bellman backup equation:
$$q_{\pi}(s,a) = r(s,a) + \sum_{a',s'} T(s' \mid s, a') \left(q_{\pi}(s',a') - \log(\pi(a' \mid s')) \right)$$

Bellman backup update operator:

$$Q(s_t, a_t) \leftarrow r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1}, a_{t+1}}[Q(s_{t+1}, a_{t+1} | s_{t+1})]$$

Soft Bellman backup update operator:

$$Q(s_t, a_t) \leftarrow r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1}, a_{t+1}} \left[Q(s_{t+1}, a_{t+1}) - \log(\pi(a_{t+1} | s_{t+1})) \right]$$

Soft Bellman backup update operator is a contraction

$$Q(s_t, a_t) \leftarrow r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1}, a_{t+1}} \left[Q(s_{t+1}, a_{t+1}) - \log(\pi(a_{t+1} | s_{t+1})) \right]$$

$$Q(s_{t}, a_{t}) \leftarrow r(s_{t}, a_{t}) + \gamma \mathbb{E}_{s_{t+1} \sim \rho} [\mathbb{E}_{a_{t+1} \sim \pi} [Q(s_{t+1}, a_{t+1}) - \log(\pi(a_{t+1} | s_{t+1}))]] \\ \leftarrow r(s_{t}, a_{t}) + \gamma \mathbb{E}_{s_{t+1} \sim \rho, a_{t+1} \sim \pi} Q(s_{t+1}, a_{t+1}) + \gamma \mathbb{E}_{s_{t+1} \sim \rho} \mathbb{E}_{a_{t+1} \sim \pi} [-\log \pi(a_{t+1} | s_{t+1})] \\ \leftarrow r(s_{t}, a_{t}) + \gamma \mathbb{E}_{s_{t+1} \sim \rho, a_{t+1} \sim \pi} Q(s_{t+1}, a_{t+1}) + \gamma \mathbb{E}_{s_{t+1} \sim \rho} \mathbb{H}(\pi(\cdot | s_{t+1}))$$

Rewrite the reward as:

$$r_{soft}(s_t, a_t) = r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \rho} \mathbf{H}(\pi(\cdot \mid s_{t+1}))$$

Then we get the old Bellman operator, which we know is a contraction

Soft Bellman backup update operator

$$Q(s_t, a_t) \leftarrow r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1}, a_{t+1}} \left[Q(s_{t+1}, a_{t+1}) - \alpha \log \pi(a_{t+1} | s_{t+1}) \right]$$

$$Q(s_{t}, a_{t}) \leftarrow r(s_{t}, a_{t}) + \gamma \mathbb{E}_{s_{t+1} \sim \rho} [\mathbb{E}_{a_{t+1} \sim \pi} [Q(s_{t+1}, a_{t+1}) - \alpha \log \pi(a_{t+1} | s_{t+1})]]$$

$$\leftarrow r(s_{t}, a_{t}) + \gamma \mathbb{E}_{s_{t+1} \sim \rho, a_{t+1} \sim \pi} Q(s_{t+1}, a_{t+1}) + \gamma \alpha \mathbb{E}_{s_{t+1} \sim \rho} \mathbb{E}_{a_{t+1} \sim \pi} [-\log \pi(a_{t+1} | s_{t+1})]$$

$$\leftarrow r(s_{t}, a_{t}) + \gamma \mathbb{E}_{s_{t+1} \sim \rho, a_{t+1} \sim \pi} Q(s_{t+1}, a_{t+1}) + \gamma \alpha \mathbb{E}_{s_{t+1} \sim \rho} H(\pi(\cdot | s_{t+1}))$$

We know that:

$$Q(s_t, a_t) \leftarrow r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim \rho}[V(s_{t+1})]$$

Which means that:

$$V(s_t) = \mathbb{E}_{a_t \sim \pi}[Q(s_t, a_t) - \alpha \log \pi(a_t | s_t)]$$

Policy iteration iterates between two steps:

1. Policy evaluation: Fix policy, apply Bellman backup operator till convergence

$$q_{\pi}(s, a) \leftarrow r(s, a) + \gamma \mathbb{E}_{s', a'} q_{\pi}(s', a')$$

4. Policy improvement: Update the policy

$$\pi'(s) \doteq \operatorname*{argmax}_{a} q_{\pi}(s, a)$$

Soft policy iteration iterates between two steps:

1. Soft policy evaluation: Fix policy, apply Bellman backup operator till convergence

$$q_{\pi}(s,a) = r(s,a) + \mathbb{E}_{s',a'} \left(q_{\pi}(s',a') - \alpha \log(\pi(a' | s')) \right)$$

This converges to q_{π}

2. Soft policy improvement: Update the policy:

$$\pi' = \arg\min_{\pi_k \in \Pi} D_{KL} \left(\pi_k(\cdot | s_t) | | \frac{\exp(Q^{\pi}(s_t, \cdot))}{Z^{\pi}(s_t)} \right)$$

Leads to a sequence of policies with monotonically increasing soft q values This so far concerns tabular methods. Next we will use function approximations for policy and action values

Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor

Let π, π' be any pair of deterministic policies such that, for all $s \in S$: $q_{\pi}(s, \pi'(s)) \ge v_{\pi}(s)$.

Then π' must be as good as or better than π , that is:

 $\mathcal{V}_{\pi'(s)\geq \mathcal{V}_{\pi}(s)}$

Review: Policy Improvement theorem for deterministic policies

Let π, π' be any pair of deterministic policies such that, for all $s \in S$: $q_{\pi}(s, \pi'(s)) \ge v_{\pi}(s)$.

Then π' must be as good as or better than π , that is:

 $\mathcal{V}_{\pi'(s)\geq \mathcal{V}_{\pi}(s)}$

$$\begin{aligned} v_{\pi}(s) &\leq q_{\pi}(s, \pi'(s)) \\ &= \mathbb{E}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s, A_t = \pi'(s)] \qquad (by (4.6)) \\ &= \mathbb{E}_{\pi'}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s] \\ &\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma q_{\pi}(S_{t+1}, \pi'(S_{t+1})) \mid S_t = s] \qquad (by (4.7)) \\ &= \mathbb{E}_{\pi'}[R_{t+1} + \gamma \mathbb{E}_{\pi'}[R_{t+2} + \gamma v_{\pi}(S_{t+2})|S_{t+1}, A_{t+1} = \pi'(S_{t+1})] \mid S_t = s] \\ &= \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 v_{\pi}(S_{t+2}) \mid S_t = s] \\ &\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 v_{\pi}(S_{t+3}) \mid S_t = s] \\ &\vdots \\ &\leq \mathbb{E}_{\pi'}[R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \cdots \mid S_t = s] \\ &= v_{\pi'}(s). \end{aligned}$$

Review: Policy Improvement theorem for deterministic policies

Let π, π' be any pair of deterministic policies such that, for all $s \in S$: $q_{\pi}(s, \pi'(s)) \ge v_{\pi}(s)$.

Then π' must be as good as or better than π , that is:

 $\mathcal{V}_{\pi'(s)\geq \mathcal{V}_{\pi}(s)}$

$$\pi'(s) \doteq \operatorname{argmax}_{a} q_{\pi}(s, a)$$

= $\operatorname{argmax}_{a} \mathbb{E}[R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s, A_t = a]$
= $\operatorname{argmax}_{a} \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_{\pi}(s')],$

$$\pi' = \arg\min_{\pi_k \in \Pi} D_{KL} \left(\pi_k(\cdot \mid s_t) \mid \left| \frac{\exp(Q^{\pi}(s_t, \cdot))}{Z^{\pi}(s_t)} \right) \right)$$

SoftMax



Use function approximations for policy and action value functions:

 $\pi_{\phi}(a_t \,|\, s_t) \qquad Q_{\theta}(s_t)$

Use function approximations for policy and action value functions: $\pi_{\phi}(a_t | s_t) \qquad Q_{\theta}(s_t)$

1. Learning the state-action value function:

$$J_Q(\theta) = \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t) \sim \mathcal{D}} \left[\frac{1}{2} \left(Q_\theta(\mathbf{s}_t, \mathbf{a}_t) - \hat{Q}(\mathbf{s}_t, \mathbf{a}_t) \right)^2 \right]$$

Semi-gradient method:

 $\hat{\nabla}_{\theta} J_Q(\theta) = \nabla_{\theta} Q_{\theta}(\mathbf{a}_t, \mathbf{s}_t) \left(Q_{\theta}(\mathbf{s}_t, \mathbf{a}_t) - \left(r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \left(Q_{\bar{\theta}}(\mathbf{s}_{t+1}, \mathbf{a}_{t+1}) - \alpha \log \left(\pi_{\phi}(\mathbf{a}_{t+1} | \mathbf{s}_{t+1}) \right) \right) \right)$

Use function approximations for policy and action value functions: $\pi_{\phi}(a_t | s_t) \quad Q_{\theta}(s_t)$

3. Learning the policy:

$$J_{\pi}(\phi) = \mathbb{E}_{\mathbf{s}_{t} \sim \mathcal{D}} \left[\mathbb{D}_{\mathrm{KL}} \left(\pi_{\phi}(\cdot | \mathbf{s}_{t}) \| \frac{\exp\left(Q_{\theta}(\mathbf{s}_{t}, \cdot)\right)}{Z_{\theta}(\mathbf{s}_{t})} \right) \right] \qquad Z_{\theta}(s_{t}) = \int_{\mathcal{A}} \exp(Q_{\theta}(s_{t}, a_{t})) da_{t}$$

$$\nabla_{\phi} J_{\pi}(\phi) = \nabla_{\phi} \mathbb{E}_{s_{t} \in D} \mathbb{E}_{a_{t} \sim \pi_{\phi}(a|s_{t})} \log \frac{\pi_{\phi}(a_{t}|s_{t})}{\frac{\exp(Q_{\theta}(s_{t}, a_{t}))}{Z_{\theta}(s_{t})}} \qquad \text{independent of \phi}$$

$$\nabla_{\phi} J_{\pi}(\phi) = \nabla_{\phi} \mathbb{E}_{s_{t} \in D, \epsilon \sim \mathcal{N}(0, I)} \log \frac{\pi_{\phi}(a_{t}|s_{t})}{\exp(Q_{\theta}(s_{t}, a_{t}))}$$

The variable w.r.t. which we take gradient parametrizes the distribution inside the expectation.

Use function approximations for policy and action value functions: $\pi_{\phi}(a_t | s_t) \quad Q_{\theta}(s_t)$

3. Learning the policy:

$$\nabla_{\phi} J_{\pi}(\phi) = \nabla_{\phi} \mathbb{E}_{s_t \in D} \mathbb{E}_{a_t \sim \pi_{\phi}(a|s_t)} \log \frac{\pi_{\phi}(a_t|s_t)}{\exp(Q_{\theta}(s_t, a_t))}$$

Reparametrization trick. The policy becomes a deterministic function of Gaussian random variables (fixed Gaussian distribution):

$$a_{t} = f_{\phi}(s_{t}, \epsilon) = \mu_{\phi}(s_{t}) + \epsilon \Sigma_{\phi}(s_{t}), \quad \epsilon \sim \mathcal{N}(0, I)$$

$$\checkmark$$

$$\nabla_{\phi} J_{\pi}(\phi) = \nabla_{\phi} \mathbb{E}_{s_{t} \in D, \epsilon \sim \mathcal{N}(0, I)} \log \frac{\pi_{\phi}(a_{t} \mid s_{t})}{\exp(Q_{\theta}(s_{t}, a_{t}))}$$





Composability of Maximum Entropy Policies

Imagine we want to satisfy two objectives at the same time, e.g., pick an object up while avoiding an obstacle. We would learn a policy to maximize the addition of the the corresponding reward functions:

$$r^{C}(s, a) = \frac{1}{C} \sum_{i=1}^{C} r_{i}(s, a)$$

MaxEnt policies permit to obtain the resulting policy's optimal Q by simply adding the constituent Qs:

$$Q_C^*(s,a) \approx \frac{1}{C} \sum_{i=1}^C Q_i^*(s,a)$$

We can theoretically bound the suboptimality of the resulting policy w.r.t. the policy trained under the addition of rewards. We cannot do this for deterministic policies.

Composable Deep Reinforcement Learning for Robotic Manipulation, Haarnoja et al.

Composable Deep Reinforcement Learning for Robotic Manipulation

Tuomas Haarnoja, Vitchyr Pong, Aurick Zhou, Murtaza Dalal, Pieter Abbeel, and Sergey Levine

> Berkeley Artificial Intelligence Research UC Berkeley