



Maximum Entropy Reinforcement Learning

Guo-Yu Yang

2020.10.26





Outline

Background

Redefinition: Bellman function

Some algorithms

Summary

Reference



Background: entropy regularized policy optimization problem

- Given K actions and the corresponding reward vector $\mathbf{q} \in \mathbb{E}^K$, and π is a probability distribution in action space. The entropy regularized policy optimization problem:

$$\max_{\pi} \{ \pi \mathbf{q} + \tau \mathcal{H}(\pi) \}$$

$\tau \geq 0$ controls the degree of exploration.

- Entropy of policy:

$$\mathcal{H}(\pi) = - \sum_{\pi_a \in \pi} \pi_a \log(\pi_a)$$

The entropy of deterministic policy is relatively low, and the entropy of random policy is relatively high.





Background: Softmax function

- The Softmax \mathcal{F}_τ function:

$$\mathcal{F}_\tau(\mathbf{q}) = \tau \log \sum_a e^{q_a/\tau}$$

- And we define the Soft-Indmax \mathbf{f}_τ function:

$$\mathbf{f}_\tau(\mathbf{q}) = \nabla \mathcal{F}_\tau(\mathbf{q}) = \frac{e^{\mathbf{q}/\tau}}{\sum_a e^{q_a/\tau}} = e^{(\mathbf{q} - \mathcal{F}_\tau(\mathbf{q}))/\tau}$$

- Soft-Indmax function gets the confidence of each action (different from Hardmax function), and it's possible to explore.





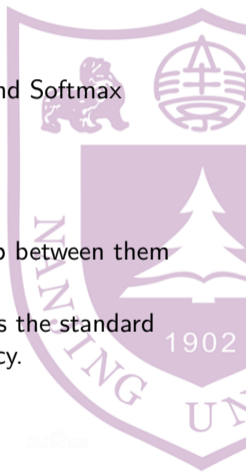
Background: the connection

- The connection between entropy regularized policy optimization problem and Softmax function:

$$\mathcal{F}_\tau(\mathbf{q}) = \max_{\pi} \{ \pi \mathbf{q} + \tau \mathcal{H}(\pi) \} = \mathbf{f}_\tau(\mathbf{q}) \cdot \mathbf{q} + \tau \mathcal{H}(\mathbf{f}_\tau(\mathbf{q}))$$

where $\pi^* = \mathbf{f}_\tau(\mathbf{q})$ and $\mathcal{F}_\tau(\mathbf{q}) = q_a - \tau \log \pi_a, \forall a$.

- The Softmax value is the upper bound on the maximum value, and the gap between them is the entropy of the policy.
- When $\tau \rightarrow 0$, The entropy regularized policy optimization problem becomes the standard expected reward objective, where the optimal solution is the hard-max policy.





Background: the connection

- Proof:

- The first equation: Let \mathcal{F}_τ^* denotes as conjugate of \mathcal{F}_τ :

$$\mathcal{F}_\tau^*(\mathbf{p}) = \sup_{\mathbf{q}} \{\mathbf{p} \cdot \mathbf{q} - \mathcal{F}_\tau(\mathbf{q})\} = \tau \mathbf{p} \log \mathbf{p}$$

For $\sum_{p \in \mathbf{p}} p = 1$. Since \mathcal{F}_τ is closed and convex, $\mathcal{F}_\tau = \mathcal{F}_\tau^{**}$:

$$\mathcal{F}_\tau(\mathbf{q}) = \sup_{\mathbf{p}} \{\mathbf{p} \cdot \mathbf{q} - \tau \mathbf{p} \log \mathbf{p}\}$$

- The second equation uses Lagrange multiplier method:

$$L = \pi(\mathbf{q} - \tau \log \pi) + \lambda(1 - \mathbf{1} \cdot \pi)$$

KKT condition:

$$1 - \mathbf{1} \cdot \pi = 0 ; \tau \log \pi = \mathbf{q} - v \quad (v = \lambda + \tau)$$





Redefinition: Bellman function

- Standard reinforcement learning objective:

$$\pi_{std}^* = \arg \max_{\pi} \sum_{t=0}^{\infty} \mathbb{E}_{\rho_{\pi}} [\gamma^t r_t]$$

- Maximum entropy reinforcement learning objective:

$$\pi_{MaxEnt}^* = \arg \max_{\pi} \sum_{t=0}^{\infty} \mathbb{E}_{\rho_{\pi}} [\gamma^t (r_t + \tau \mathcal{H}(\pi))]$$

- Original Q-function:

$$Q^*(s_t, a_t) = r_t + \mathbb{E}_{\rho_{\pi^*}} \left[\sum_{l=1}^{\infty} \gamma^l r_{t+l} \right]$$

- Soft Q-function:

$$Q_{soft}^*(s_t, a_t) = r_t + \mathbb{E}_{\rho_{\pi^*}} \left[\sum_{l=1}^{\infty} \gamma^l (r_{t+l} + \tau \mathcal{H}(\pi^*)) \right]$$



Redefinition: Bellman function

- Original optimal value function:

$$V^*(s) = \max_{a \in \mathcal{A}} Q^*(s, a)$$

- Soft-value function:

$$V_{soft}^*(s) = \tau \log \int_{\mathcal{A}} \exp\left(\frac{1}{\tau} Q_{soft}^*(s, \mathbf{a}')\right) d\mathbf{a}'$$

- According to the above re-definition, we have following conclusions:

- 1 Maximum entropy policy given by (Boltzmann distribution):

$$\pi_{MaxEnt}^*(a_t | s_t) = \exp\left(\frac{1}{\tau} (Q_{soft}^*(s_t, a_t) - V_{soft}^*(s_t))\right)$$

- 2 The soft Bellman Equation:

$$Q_{soft}^*(s_t, a_t) = r(s_t, a_t) + \gamma \mathbb{E}_{\pi^*} [V_{soft}^*(s_{t+1})]$$



Redefinition: policy improvement theorem

- Policy improvement theorem: Given a policy π , define a new policy $\tilde{\pi}$ as

$$\tilde{\pi}(\cdot|s) \propto \exp(Q_{soft}^{\pi}(s, \cdot))$$

Assume that Q is bounded and $\int \exp(Q(s, a)) da$ is bounded for any s . Then $Q^{\tilde{\pi}}(s, a) \geq Q^{\pi}(s, a), \forall s, a$.

- Proof:
 - Now, we proof the next inequality:

$$\tau \mathcal{H}(\pi(\cdot|s)) + \mathbb{E}_{a \sim \pi}[Q_{soft}^{\pi}(s, a)] \leq \tau \mathcal{H}(\tilde{\pi}(\cdot|s)) + \mathbb{E}_{a \sim \tilde{\pi}}[Q_{soft}^{\pi}(s, a)]$$

According to the principle of normalization:

$$\begin{aligned} \tilde{\pi}(\cdot|s) &= \exp\left(\frac{1}{\tau}(Q_{soft}^{\pi}(s, \cdot) - V_{soft}^{\pi}(s))\right) \\ Q_{soft}^{\pi}(s, \cdot) &= V_{soft}^{\pi}(s) + \tau \log(\tilde{\pi}(\cdot|s)) \end{aligned}$$

Tuomas Haarnoja et al. "Reinforcement Learning with Deep Energy-Based Policies". In: Proceedings of the 34th International Conference on Machine Learning, ICML 2017.



Redefinition: policy improvement theorem

- Proof:

- Therefore, we have the following equation:

$$\begin{aligned}
 \tau \mathcal{H}(\pi(\cdot|s)) + \mathbb{E}_{a \sim \pi}[Q_{soft}^\pi(s, a)] &= \tau \mathcal{H}(\pi(\cdot|s)) + V_{soft}^\pi(s) + \tau \pi(\cdot|s) \log(\tilde{\pi}(\cdot|s)) \\
 &= -\tau D_{KL}(\pi(\cdot|s) \parallel \tilde{\pi}(\cdot|s)) + V_{soft}^\pi(s) \\
 &\leq \tau \mathcal{H}(\tilde{\pi}(\cdot|s)) + V_{soft}^\pi(s) + \tau \tilde{\pi}(\cdot|s) \log(\tilde{\pi}(\cdot|s)) \\
 &= \tau \mathcal{H}(\tilde{\pi}(\cdot|s)) + \mathbb{E}_{a \sim \tilde{\pi}}[Q_{soft}^\pi(s, a)]
 \end{aligned}$$

- Then we can show that:

$$\begin{aligned}
 Q_{soft}^\pi(s, a) &= \mathbb{E}_{s_1}[r_0 + \gamma(\mathcal{H}(\pi(\cdot|s_1)) + \mathbb{E}_{a_1 \sim \pi}[Q_{soft}^\pi(s_1, a_1)])] \\
 &\leq \mathbb{E}_{s_1}[r_0 + \gamma(\mathcal{H}(\tilde{\pi}(\cdot|s_1)) + \mathbb{E}_{a_1 \sim \tilde{\pi}}[Q_{soft}^\pi(s_1, a_1)])] \\
 &= \mathbb{E}_{s_1}[r_0 + \gamma(\mathcal{H}(\tilde{\pi}(\cdot|s_1)) + r_1)] + \gamma^2 \mathbb{E}_{s_2}[\mathcal{H}(\pi(\cdot|s_2)) + \mathbb{E}_{a_2 \sim \pi}[Q_{soft}^\pi(s_2, a_2)]] \\
 &\dots \\
 &\leq \mathbb{E}_{\rho \sim \tilde{\pi}}[r_0 + \sum_{t=1} \gamma^t (\mathcal{H}(\tilde{\pi}(\cdot|s_t)) + r_t)] \\
 &= Q^{\tilde{\pi}}(s, a)
 \end{aligned}$$





Redefinition: Bellman function

- According to policy improvement theorem, after several iterations we can get:

$$\pi_{MaxEnt}^*(a_t|s_t) = \exp\left(\frac{1}{\tau}(Q_{soft}^*(s_t, a_t) - V_{soft}^*(s_t))\right)$$

- Proof of soft Bellman Equation:

- Soft Q-function rewritten as:

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim p}[\tau \mathcal{H}(\pi(\cdot|s')) + \mathbb{E}_{a' \sim \pi(\cdot|s')}(Q^\pi(s', a'))]$$

- according to policy improvement theorem:

$$Q^\pi(s, a) = V^\pi(s) + \tau \log(\pi(a|s))$$

- Combine the two equation:

$$Q^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim p}[V^\pi(s')]$$





Algorithm: Soft Q-Iteration

- Let Q and V be bounded. The fixed-point iteration:

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

$$V(s_t) \leftarrow \tau \log \left(\int_{\mathcal{A}} \exp\left(\frac{1}{\tau} Q(s_t, a')\right) da' \right), \forall s_t$$

converges to Q^* and V^* respectively.

- Disadvantage:
 - The soft Bellman backup cannot be performed exactly in continuous or large state and action spaces.
 - Sampling from the energy-based model in $\tilde{\pi}(\cdot|s) \propto \exp(Q_{soft}^{\pi}(s, \cdot))$ is intractable in general.



Algorithm: Soft Q-Iteration

- Proof:

- We show that the soft value iteration operator \mathcal{T} , defined as:

$$\mathcal{T}Q(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim p_s} [\log \int \exp Q(s', a') da']$$

Define a norm $\|Q_1 - Q_2\| = \max_{s,a} |Q_1(s, a) - Q_2(s, a)|$. Suppose $\varepsilon = \|Q_1 - Q_2\|$:

$$\begin{aligned} \log \int \exp(Q_1(s', a')) da' &\leq \log \int \exp(Q_2(s', a') + \varepsilon) da' \\ &= \log(\exp(\varepsilon) \int \exp Q_2(s', a') da') \\ &= \varepsilon + \log \int \exp(Q_2(s', a')) da' \end{aligned}$$

Therefore, $\|\mathcal{T}Q_1 - \mathcal{T}Q_2\| \leq \gamma\varepsilon = \gamma\|Q_1 - Q_2\|$. So \mathcal{T} is a contraction.



Algorithm: Soft Q-learning

In order to solve the above problem of Soft Q-iteration, we use stochastic optimization problem to model. The following is the pseudocode of Soft Q-learning:

Algorithm 1 Soft Q-learning

Input: θ, ϕ

- 1: $\bar{\theta} \leftarrow \theta, \bar{\phi} \leftarrow \phi$
- 2: $\mathcal{D} \leftarrow \emptyset$
- 3: **for** each epoch **do**
- 4: **for** each t **do**
- 5: $a_t = f^\phi(\xi, s_t)$, where $\xi \sim \mathcal{N}(0, I)$
- 6: $s_{t+1} \sim p(s_{t+1} | s_t, a_t)$
- 7: $\mathcal{D} \leftarrow \mathcal{D} \cup \{(s_t, a_t, r(s_t, a_t), s_{t+1})\}$
- 8: **Update soft Q-function parameters**
- 9: **Update policy**
- 10: **end for**
- 11: **if** epoch % update_interval = 0 **then**
- 12: $\bar{\theta} \leftarrow \theta, \bar{\phi} \leftarrow \phi$
- 13: **end if**
- 14: **end for**
- 15: **return** θ, ϕ

Tuomas Haarnoja et al. "Reinforcement Learning with Deep Energy-Based Policies" . In: Proceedings of the 34th International Conference on Machine Learning, ICML 2017.

Algorithm: Soft Q-learning

- Firstly, we model the soft Q-function with a function approximator with parameters θ , $Q_\theta(s_t, a_t)$. We will optimize the soft Bellman error $|\mathcal{T}Q - Q|$ to find optimal Q-function.
- Update soft Q-learning parameters:

- ① Sample a minibatch from the replay memory: $\{s_t^{(i)}, a_t^{(i)}, r_t^{(i)}, s_{t+1}^{(i)}\}_{i=0}^N \sim \mathcal{D}$
- ② Sample M uniform actions: $\{a^{(i,j)}\}_{j=0}^M \sim q_{a'}$
- ③ Compute empirical soft value with next equation:

$$V_{\bar{\theta}}(s_{t+1}^{(i)}) = \tau \log \mathbb{E}_{q_{a'}} \left[\frac{\exp(\frac{1}{\tau} Q_\theta(s_t, a'))}{q_{a'}(a')} \right]$$

It use importance sample, where $q_{a'}$ is a non-zero distribution over the action space.

- ④ Compute empirical gradient of next equation (equivalent to Sort Q-iteration):

$$g_\theta^{(i)} \leftarrow \nabla_\theta \left(\mathbb{E}_{s_t \sim q_s, a_t \sim q_a} \left[\frac{1}{2} (Q_\theta(s_t, a_t) - (r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1} \sim p} [V_{\bar{\theta}}(s_{t+1})]))^2 \right] \right)$$

- ⑤ Update θ : $\theta \leftarrow ADAM(\theta, g_\theta^{(i)})$

Algorithm: Soft Q-learning

- Secondly, we want to learn a state-conditioned stochastic neural network $a_t = f^\phi(\xi, s_t)$.
- Update policy process:

- Sample M latents for each state: $\{\xi^{(i,j)}\}_{j=0}^M \sim \mathcal{N}(0, I)$
- Evaluate actions: $a_t = f^\phi(\xi^{(i,j)}, s_t^{(i)})$
- Evaluate empirical Stein variational gradient using next equation:

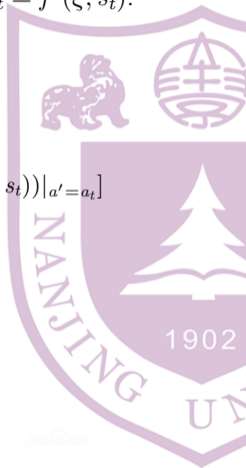
$$\Delta f^\phi(\cdot; s_t) = \mathbb{E}_{a_t \sim \pi_\phi} [\kappa(a_t, f_\phi(\cdot; s_t)) \nabla_{a'} Q_\theta(s_t, a')|_{a'=a_t} + \tau \nabla_{a'} \kappa(a', f^\phi(\cdot; s_t))|_{a'=a_t}]$$

- Compute empirical estimate $g_\phi^{(i)} = \hat{\nabla}_\phi J_\pi$ of next equation:

$$\frac{\partial J_\pi(\phi; s_t)}{\partial \phi} \propto \mathbb{E}_\xi \left[\Delta f^\phi(\xi; s_t) \frac{\partial f^\phi(\xi, s_t)}{\partial \phi} \right]$$

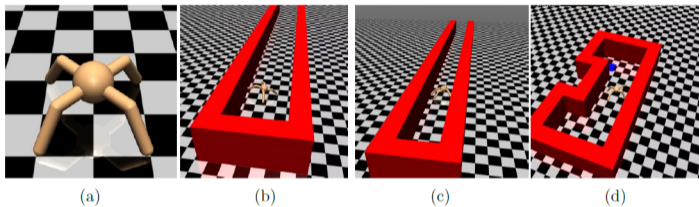
$$J_\pi(\phi; s_t) = D_{KL} \left(\pi_\phi(\cdot | s_t) \parallel \exp\left(\frac{1}{\tau} (Q_\theta(s_t, \cdot) - V_\theta(s_t))\right) \right)$$

- Update ϕ : $\phi \leftarrow \text{ADAM}(\phi, g_\phi^{(i)})$

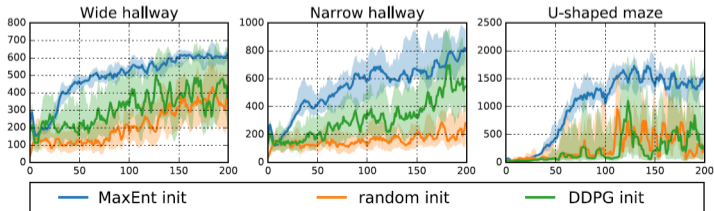


Algorithm: Soft Q-learning

- Experimental environment:



- Experimental environment:



Algorithm: Maximum entropy monte-carlo planning

- UCT can only guarantee a polynomial convergence rate of finding the best action at the root in MCTS.
- MENTS (Maximum Entropy for Tree Search) improves UCT to exponential convergence speed:
 - ① Using E2W (empirical exponential weight) as the tree policy.
 - ② Evaluating each search node by softmax values back-propagated from simulations.
- The stochastic softmax bandit:
 The target is finding the policy with maximum softmax value in action set \mathcal{A} :

$$V_{soft}^* = \mathcal{F}_\tau(r)$$

An important assumption: reward function of action, $r(a)$, are from σ^2 -subgaussian.





Algorithm: Maximum entropy monte-carlo planning

- Empirical exponential weight (E2W): enforce enough exploration to guarantee good estimation of q , and make the policy converge to π^* asymptotically. The algorithm selects an action by sampling from the distribution:

$$\pi_t(a) = (1 - \lambda_t)\mathbf{f}_\tau(\mathbf{q})(a) + \lambda_t \frac{1}{|\mathcal{A}|}$$

where $\lambda_t = \epsilon|\mathcal{A}|/\log(t+1)$ is a decay rate for exploration, with exploration parameter $\epsilon > 0$.

- Theoretical guarantee of E2W: In the softmax stochastic bandit problem:

$$\lim_{t \rightarrow \infty} t\mathcal{E}^t = \frac{\sigma^2}{\tau^2} \left(\sum_a \exp(r(a)/\tau) \right)^2$$

where t is number of iterations, σ is the assumption of the subsequent distribution, \mathcal{E}_t is the mean square error in stochastic bandit problem.





Algorithm: Maximum entropy monte-carlo planning

- MENTS policy:

$$\pi_t(a|s) = (1 - \lambda_s) \mathbf{f}_\tau(\mathbf{Q}_{soft}(s))(a) + \lambda_s \frac{1}{|\mathcal{A}|}$$

where $\lambda_s = \epsilon |\mathcal{A}| / \log(\sum_a N(s, a) + 1)$. \mathbf{Q}_{soft} denote a \mathcal{A} -dimensional vector with components $Q_{soft}(s, a)$. Q-values using the softmax backup:

$$Q_{soft}(s_t, a_t) = \begin{cases} r(s_t, a_t) + R & t = T - 1 \\ r(s_t, a_t) + \mathcal{F}_\tau(\mathbf{Q}_{soft}(s_{t+1})) & t < T - 1 \end{cases}$$

R is the evaluation value by rollout policy. Finally, MENTS proposes the action with the maximum estimated softmax value.





Algorithm: Maximum entropy monte-carlo planning

- Theoretical guarantee:

Let a_t be the action returned by MENTS at iteration t . Then for large enough t with some constant C :

$$\mathbb{P}(a_t \neq a^*) \leq Ct \exp\left(-\frac{t}{(\log t)^3}\right)$$

- Experimental effect:

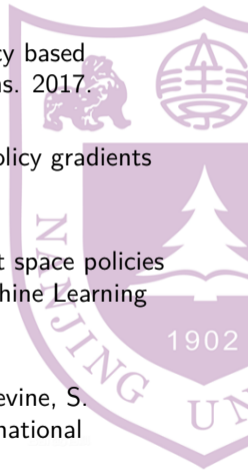
Table 1: Performance comparison of Atari games playing.

Agent	<i>BeamRider</i>	<i>Breakout</i>	<i>Q*bert</i>	<i>Seaquest</i>	<i>SpaceInvaders</i>
DQN	19280	345	14558	1142	625
UCT	21952	367	16010	1129	656
MENTS	18576	386	18336	1161	1503



Other Algorithm I

- ① Path Consistency Learning (PCL):
 Reference: Nachum, Ofir, et al. "Bridging the gap between value and policy based reinforcement learning." Advances in Neural Information Processing Systems. 2017.
- ② Prove the equivalence of soft Q-learning and policy gradients
 Reference: Schulman, J., Abbeel, P., and Chen, X. Equivalence between policy gradients and soft Q-learning. arXiv preprint arXiv:1704.06440, 2017.
- ③ Hierarchical maximum entropy reinforcement learning
 Reference: Haarnoja, T., Hartikainen, K., Abbeel, P., and Levine, S. Latent space policies for hierarchical reinforcement learning. In International Conference on Machine Learning (ICML), 2018.
- ④ Combine Soft Q-learning policies to produce better policy
 Reference: Haarnoja, T., Pong, V., Zhou, A., Dalal, M., Abbeel, P., and Levine, S. Composable deep reinforcement learning for robotic manipulation. In International Conference on Robotics and Automation (ICRA). IEEE, 2018.





Other Algorithm II

- ⑤ Soft Actor-Critic algorithm
 Reference: Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In International Conference on Machine Learning (ICML), 2018.
- ⑥ Dynamically and automatically tuning the temperature parameter in Soft Actor-Critic Algorithm
 Reference: Haarnoja, Tuomas, et al. "Soft actor-critic algorithms and applications." arXiv preprint arXiv:1812.05905. 2018.
- ⑦ Soft policy gradient method
 Reference: Shi, Wenjie et al. "Soft Policy Gradient Method for Maximum Entropy Deep Reinforcement Learning." IJCAI (2019).





Summary

- 1 Entropy regularized policy optimization problem: trade off exploration-exploitation directly
- 2 Bellman equation redefinition: combine with entropy regularized policy optimization problem and softmax function
- 3 Simple algorithm: soft Q-iteration
- 4 More practical algorithm: soft Q-learning
- 5 Maximum Entropy for Tree Search: variations of MCTS

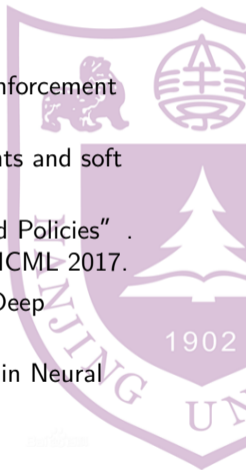




Reference

Note: Papers in this area mainly come from UC Berkeley, U alberta.

- 1 Nachum, Ofir, et al. "Bridging the gap between value and policy based reinforcement learning." Advances in Neural Information Processing Systems. 2017.
- 2 Schulman, J., Abbeel, P., and Chen, X. Equivalence between policy gradients and soft Q-learning. arXiv preprint arXiv:1704.06440, 2017.
- 3 Tuomas Haarnoja et al. "Reinforcement Learning with Deep Energy-Based Policies" . In:Proceedings of the 34th International Conference on Machine Learning, ICML 2017.
- 4 Haarnoja, Tuomas. Acquiring Diverse Robot Skills via Maximum Entropy Deep Reinforcement Learning. Diss. UC Berkeley, 2018.
- 5 Xiao, Chenjun, et al. "Maximum entropy monte-carlo planning." Advances in Neural Information Processing Systems. 2019.





Thankyou!

