

Lecture 5b: Maximum Entropy RL

CS885 Reinforcement Learning

2025-01-21

Complementary readings:

Haarnoja, Tang, Abbeel, Levine (2017) Reinforcement Learning with Deep Energy-Based Policies, ICML.

Haarnoja, Zhou, Abbeel, Levine (2018) Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor, ICML.

Pascal Poupart

David R. Cheriton School of Computer Science



Maximum Entropy RL

- Why do several implementations of important RL baselines (e.g., A2C, PPO) add an entropy regularizer?
- Why is maximizing entropy desirable in RL?
- What is the Soft Actor Critic algorithm?

Reinforcement Learning

Deterministic Policies

- There always exists an optimal deterministic policy
- Search space is smaller for deterministic than stochastic policies
- Practitioners prefer deterministic policies

Stochastic Policies

- Search space is continuous for stochastic policies (helps with gradient descent)
- **More robust (less likely to overfit)**
- **Naturally incorporate exploration**
- **Facilitate transfer learning**
- **Mitigate local optima**

Encouraging Stochasticity

Standard MDP

- States: S
- Actions: A
- Reward: $R(s, a)$
- Transition: $\Pr(s'|s, a)$
- Discount: γ

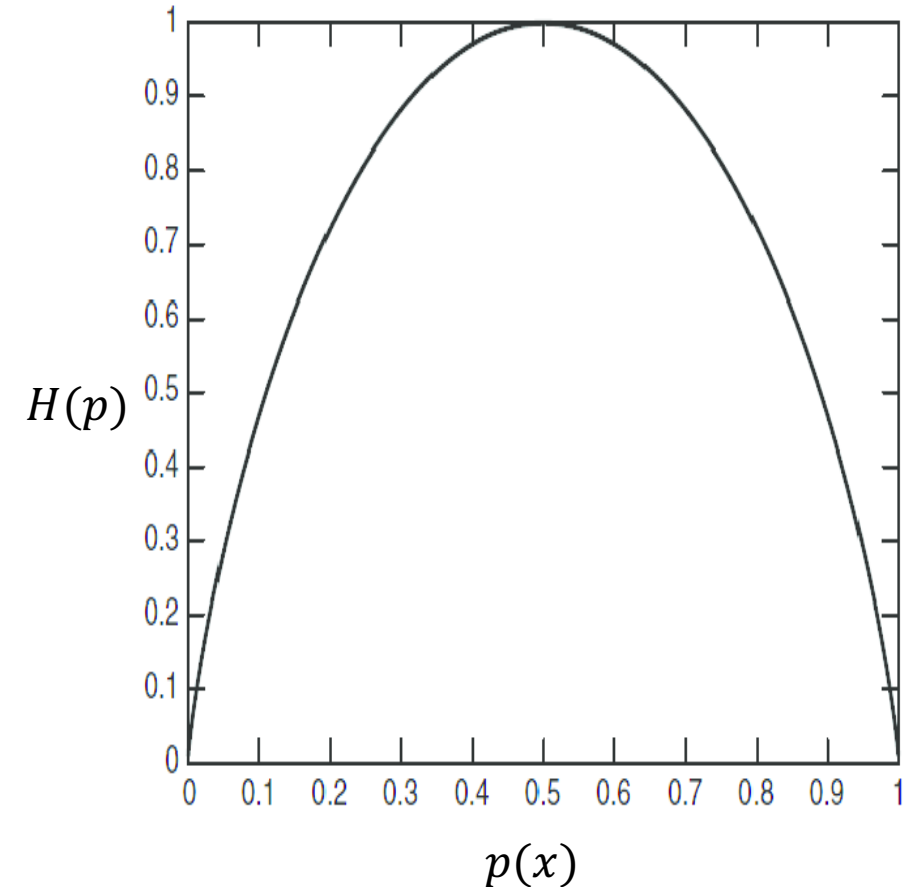
Soft MDP

- States: S
- Actions: A
- Reward: $R(s, a) + \lambda H(\pi(\cdot |s))$
- Transition: $\Pr(s'|s, a)$
- Discount: γ

Entropy

- Measure of uncertainty
 - Information theory: expected # of bits needed to communicate the result of a sample

$$H(p) = -\sum_x p(x) \log p(x)$$



Optimal Policy

▪ Standard MDP:
$$\pi^* = \operatorname{argmax}_{\pi} \sum_{n=0}^N \gamma^n E_{s_n, a_n | \pi} [R(s_n, a_n)]$$

▪ Soft MDP:
$$\pi_{soft}^* = \operatorname{argmax}_{\pi} \sum_{n=0}^N \gamma^n E_{s_n, a_n | \pi} [R(s_n, a_n) + \lambda H(\pi(\cdot | s_n))]$$




Maximum entropy policy
Entropy regularized policy

Q-function

- Standard MDP

$$Q^\pi(s_0, a_0) = R(s_0, a_0) + \sum_{n=1}^{\infty} \gamma^n E_{s_n, a_n | s_0, a_0, \pi} [R(s_n, a_n)]$$

- Soft MDP

$$Q_{soft}^\pi(s_0, a_0) = R(s_0, a_0) + \sum_{n=1}^{\infty} \gamma^n E_{s_n, a_n | s_0, a_0, \pi} [R(s_n, a_n) + \lambda H(\pi(\cdot | s_n))]$$


NB: **No entropy** with first reward term
since action is not chosen according to π

Greedy Policy

- Standard MDP (deterministic policy)

$$\pi_{greedy}(s) = \operatorname{argmax}_a Q(s, a)$$

- Soft MDP (stochastic policy)

$$\begin{aligned}\pi_{greedy}(\cdot | s) &= \operatorname{argmax}_{\pi} \sum_a \pi(a|s) Q(s, a) + \lambda H(\pi(\cdot | s)) \\ &= \frac{\exp(Q(s, \cdot) / \lambda)}{\sum_a \exp(Q(s, a) / \lambda)} = \operatorname{softmax}(Q(s, \cdot) / \lambda)\end{aligned}$$

when $\lambda \rightarrow 0$ then *softmax* becomes regular max

Derivation

- Concave objective (can find global maximum)

$$\begin{aligned} J(\pi, Q) &= \sum_a \pi(a|s)Q(s, a) + \lambda H(\pi(\cdot |s)) \\ &= \sum_a \pi(a|s)[Q(s, a) - \lambda \log \pi(a|s)] \end{aligned}$$

- Partial derivative: $\frac{\partial J}{\partial \pi(a|s)} = Q(s, a) - \lambda[\log \pi(a|s) + 1]$

- Setting the derivative to 0 and isolating $\pi(a|s)$ yields

$$\pi(a|s) = \exp(Q(s, a)/\lambda - 1) \propto \exp(Q(s, a)/\lambda)$$

- Hence $\pi_{greedy}(\cdot |s) = \frac{\exp(Q(s, \cdot)/\lambda)}{\sum_a \exp(Q(s, a)/\lambda)} = \text{softmax}(Q(s, \cdot)/\lambda)$

Greedy Value Function

- What is the value function induced by the greedy policy?
- Standard MDP: $V(s) = \max_a Q(s, a)$
- Soft MDP: $V_{soft}(s) = \lambda H \left(\pi_{greedy}(\cdot | s) \right) + \sum_a \pi_{greedy}(a | s) Q_{soft}(s, a)$

$$= \lambda \log \sum_a \exp \left(\frac{Q_{soft}(s, a)}{\lambda} \right) = \widetilde{\max}_a Q_{soft}(s, a)$$

when $\lambda \rightarrow 0$ then $\widetilde{\max}_a$ becomes regular max

Derivation

$$V_{soft}(s) = \lambda H \left(\pi_{greedy}(\cdot | s) \right) + \sum_a \pi_{greedy}(a | s) Q_{soft}(s, a)$$

$$\text{since } \pi_{greedy}(a | s) = \frac{\exp(Q_{soft}(s, a) / \lambda)}{\sum_{a'} \exp(Q_{soft}(s, a') / \lambda)}$$

$$= \lambda H \left(\pi_{greedy}(\cdot | s) \right) + \sum_a \pi_{greedy}(a | s) \lambda \left[\log \pi_{greedy}(a | s) + \log \sum_{a'} \exp \left(\frac{Q_{soft}(s, a')}{\lambda} \right) \right]$$

$$= \lambda H \left(\pi_{greedy}(\cdot | s) \right) + \lambda \sum_a \pi_{greedy}(a | s) \log \pi_{greedy}(a | s) + \lambda \log \sum_{a'} \exp \left(\frac{Q_{soft}(s, a')}{\lambda} \right)$$

$$= \lambda H \left(\pi_{greedy}(\cdot | s) \right) - \lambda H \left(\pi_{greedy}(\cdot | s) \right) + \lambda \log \sum_{a'} \exp \left(\frac{Q_{soft}(s, a')}{\lambda} \right)$$

$$= \lambda \log \sum_{a'} \exp \left(\frac{Q_{soft}(s, a')}{\lambda} \right)$$

$$= \widetilde{\max}_a \lambda Q_{soft}(s, a)$$

Soft Q-Value Iteration

SoftQValueIteration(MDP, λ)

Initialize π_0 to any policy

$i \leftarrow 0$

Repeat

$$Q_{soft}^{i+1}(s, a) \leftarrow R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) \widetilde{\max}_{a'} Q_{soft}^i(s', a')$$

$i \leftarrow i + 1$

Until $\|Q_{soft}^i(s, a) - Q_{soft}^{i-1}(s, a)\|_{\infty} \leq \epsilon$

Extract policy: $\pi_{greedy}(\cdot | s) = \text{softmax}(Q_{soft}^i(s, \cdot) / \lambda)$

Soft Bellman equation: $Q_{soft}^*(s, a) = R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) \widetilde{\max}_{a'} Q_{soft}^*(s', a')$

Soft Q-learning

- Q-learning based on Soft Q-Value Iteration
- Replace expectations by samples
- Represent Q-function by a function approximator (e.g., neural network)
- Do gradient updates based on temporal differences

Soft Q-learning (Soft Variant of DQN)

Initialize weights \mathbf{w} and $\bar{\mathbf{w}}$ at random in $[-1,1]$

Observe current state s

Loop

Select action a and execute it

Receive immediate reward r , observe new state s'

Add (s, a, s', r) to experience buffer

Sample mini-batch of experiences from buffer

For each experience $(\hat{s}, \hat{a}, \hat{s}', \hat{r})$ in mini-batch

$$\text{Gradient: } \frac{\partial \text{Err}}{\partial \mathbf{w}} = \left[Q_{\mathbf{w}}^{\text{soft}}(\hat{s}, \hat{a}) - \hat{r} - \gamma \overline{\max_{\hat{a}'} Q_{\mathbf{w}}^{\text{soft}}(\hat{s}', \hat{a}')} \right] \frac{\partial Q_{\mathbf{w}}^{\text{soft}}(\hat{s}, \hat{a})}{\partial \mathbf{w}}$$

$$\text{Update weights: } \mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial \text{Err}}{\partial \mathbf{w}}$$

Update state: $s \leftarrow s'$

Every c steps, update target: $\bar{\mathbf{w}} \leftarrow \mathbf{w}$

Soft Actor Critic

- In practice, actor critic techniques tend to perform better than Q-learning.
- Can we derive a soft actor-critic algorithm?
- Yes, idea:
 - Critic: **soft Q-function**
 - Actor: (greedy) **softmax policy**

Soft Policy Iteration

Initialize π_0 to any policy, $i \leftarrow 0$

Repeat

Policy evaluation:

Repeat until convergence

$$Q_{soft}^{\pi_i}(s, a) \leftarrow R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) \left[\sum_{a'} \pi_i(a'|s') Q_{soft}^{\pi_i}(s', a') + \lambda H(\pi_i(\cdot |s')) \right] \quad \forall s, a$$

Policy improvement:

$$\pi_{i+1}(a|s) \leftarrow \mathit{softmax} \left(Q_{soft}^{\pi_i}(s, a) / \lambda \right) = \frac{\exp(Q_{soft}^{\pi_i}(s, a) / \lambda)}{\sum_{a'} \exp(Q_{soft}^{\pi_i}(s, a') / \lambda)} \quad \forall s, a$$

$i \leftarrow i + 1$

Until $\left\| Q_{soft}^{\pi_i}(s, a) - Q_{soft}^{\pi_{i-1}}(s, a) \right\|_{\infty} \leq \epsilon$

Policy Improvement

Theorem 1: Let $Q_{soft}^{\pi_i}(s, a)$ be the Q-function of π_i
Let $\pi_{i+1}(a|s) = \text{softmax}\left(Q_{soft}^{\pi_i}(s, a)/\lambda\right)$
Then $Q_{soft}^{\pi_{i+1}}(s, a) \geq Q_{soft}^{\pi_i}(s, a) \forall s, a$

Proof: first show that

$$\sum_a \pi_i(a|s) Q_{soft}^{\pi_i}(s, a) + \lambda H(\pi_i(\cdot | s)) \leq \sum_a \pi_{i+1}(a|s) Q_{soft}^{\pi_i}(s, a) + \lambda H(\pi_{i+1}(\cdot | s))$$

then use this inequality to show that

$$Q_{soft}^{\pi_{i+1}}(s, a) \geq Q_{soft}^{\pi_i}(s, a) \forall s, a$$

Inequality Derivation

$$\begin{aligned} & \sum_a \pi_i(a|s) Q_{soft}^{\pi_i}(s, a) + \lambda H(\pi_i(\cdot | s)) \\ &= \sum_a \pi_i(a|s) \left[Q_{soft}^{\pi_i}(s, a) - \lambda \log \pi_i(a|s) \right] \quad \text{since } \pi_{i+1}(a|s) = \frac{\exp(Q_{soft}^{\pi_i}(s, a)/\lambda)}{\sum_{a'} \exp(Q_{soft}^{\pi_i}(s, a')/\lambda)} \\ &= \sum_a \pi_i(a|s) \left[\lambda \log \pi_{i+1}(a|s) - \lambda \log \sum_{a'} \exp(Q_{soft}^{\pi_i}(s, a')/\lambda) - \lambda \log \pi_i(a|s) \right] \\ &= \lambda \sum_a \pi_i(a|s) \left[\log \frac{\pi_{i+1}(a|s)}{\pi_i(a|s)} + \log \sum_{a'} \exp(Q_{soft}^{\pi_i}(s, a')/\lambda) \right] \\ &= -\lambda KL(\pi_{i+1} || \pi_i) + \lambda \sum_a \pi_i(a|s) \log \sum_{a'} \exp(Q_{soft}^{\pi_i}(s, a')/\lambda) \\ &\leq \lambda \sum_a \pi_i(a|s) \log \sum_{a'} \exp(Q_{soft}^{\pi_i}(s, a')/\lambda) \\ &= \sum_a \pi_{i+1}(a|s) \lambda \log \sum_{a'} \exp(Q_{soft}^{\pi_i}(s, a')/\lambda) \quad \text{since } \pi_{i+1}(a|s) = \frac{\exp(Q_{soft}^{\pi_i}(s, a)/\lambda)}{\sum_{a'} \exp(Q_{soft}^{\pi_i}(s, a')/\lambda)} \\ &= \sum_a \pi_{i+1}(a|s) \left[Q_{soft}^{\pi_i}(s, a) - \lambda \log \pi_{i+1}(s, a) \right] \\ &= \sum_a \pi_{i+1}(a|s) Q_{soft}^{\pi_i}(s, a) + \lambda H(\pi_{i+1}(\cdot | s)) \end{aligned}$$

Proof Derivation

$$Q_{soft}^{\pi_i}(s, a) = R(s, a) + \gamma E_{s'} \left[E_{a' \sim \pi_i} \left[Q_{soft}^{\pi_i}(s', a') \right] + \lambda H(\pi_i(\cdot | s')) \right]$$

$$\text{since } E_{a' \sim \pi_i} \left[Q_{soft}^{\pi_i}(s', a') \right] + \lambda H(\pi_i(\cdot | s')) \leq E_{a' \sim \pi_{i+1}} \left[Q_{soft}^{\pi_i}(s', a') \right] + \lambda H(\pi_{i+1}(\cdot | s'))$$

$$\leq R(s, a) + \gamma E_{s'} \left[E_{a' \sim \pi_{i+1}} \left[Q_{soft}^{\pi_i}(s', a') \right] + \lambda H(\pi_{i+1}(\cdot | s')) \right]$$

$$\leq \dots \quad \text{repeatedly apply}$$

$$\leq \dots \quad Q_{soft}^{\pi_i}(s', a') \leq R(s', a') + \gamma E_{s''} \left[E_{a'' \sim \pi_{i+1}} \left[Q_{soft}^{\pi_i}(s'', a'') \right] + \lambda H(\pi_{i+1}(\cdot | s'')) \right]$$

$$\leq Q_{soft}^{\pi_{i+1}}(s, a)$$

Convergence to Optimal Q_{soft}^* and π_{soft}^*

- Theorem 2: When $\epsilon = 0$,
soft policy iteration converges to optimal Q_{soft}^* and π_{soft}^* .
- Proof:
 - We know that $Q^{\pi_{i+1}}(s, a) \geq Q^{\pi_i}(s, a) \forall s, a$ according to Theorem 1
 - Since the Q-functions are upper bounded by $\left(\max_{s,a} R(s, a) + H(\text{uniform})\right) / (1 - \gamma)$
then soft policy iteration converges
 - At convergence, $Q^{\pi_{i-1}} = Q^{\pi_i}$ and therefore the Q-function satisfies Bellman's equation:

$$Q_{soft}^{\pi_{i-1}}(s, a) = Q_{soft}^{\pi_i}(s, a) = R(s, a) + \gamma \sum_{s'} \Pr(s'|s, a) \max_{a'} Q_{soft}^{\pi_{i-1}}(s', a')$$

Soft Actor-Critic

- RL version of soft policy iteration
- Use neural networks to represent policy and value functions
- At each policy improvement step, project new policy in the space of parameterized neural nets

Soft Actor-Critic (SAC)

Initialize weights \mathbf{w} , $\bar{\mathbf{w}}$, θ at random in $[-1,1]$

Observe current state s

Loop

Sample action $a \sim \pi_\theta(\cdot | s)$ and execute it

Receive immediate reward r , observe new state s'

Add (s, a, s', r) to experience buffer

Sample mini-batch of experiences from buffer

For each experience $(\hat{s}, \hat{a}, \hat{s}', \hat{r})$ in mini-batch

Sample $\hat{a}' \sim \pi_\theta(\cdot | \hat{s}')$

$$\text{Gradient: } \frac{\partial \text{Err}}{\partial \mathbf{w}} = \left[Q_{\mathbf{w}}^{\text{soft}}(\hat{s}, \hat{a}) - \hat{r} - \gamma [Q_{\bar{\mathbf{w}}}^{\text{soft}}(\hat{s}', \hat{a}') + \lambda H(\pi_\theta(\cdot | \hat{s}'))] \right] \frac{\partial Q_{\mathbf{w}}^{\text{soft}}(\hat{s}, \hat{a})}{\partial \mathbf{w}}$$

$$\text{Update weights: } \mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial \text{Err}}{\partial \mathbf{w}}$$

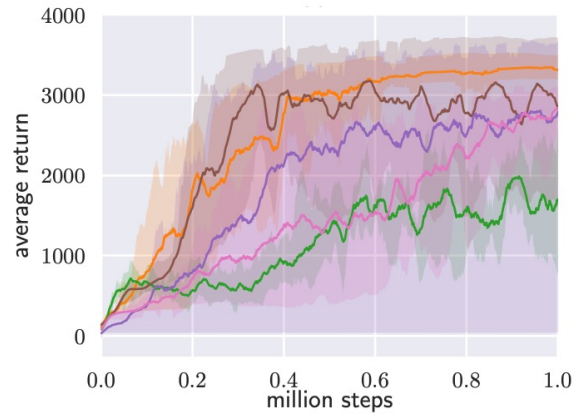
$$\text{Update policy: } \theta \leftarrow \theta - \alpha \frac{\partial_{\text{KL}}(\pi_\theta | \text{softmax}(Q_{\bar{\mathbf{w}}}^{\text{soft}} / \lambda))}{\partial \theta}$$

Update state: $s \leftarrow s'$

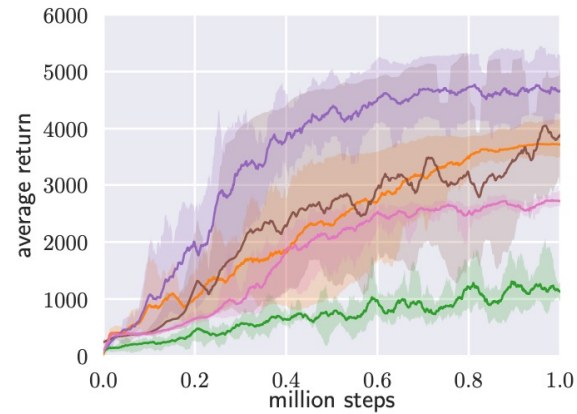
Every c steps, update target: $\bar{\mathbf{w}} \leftarrow \mathbf{w}$

Empirical Results

Comparison on several robotics tasks

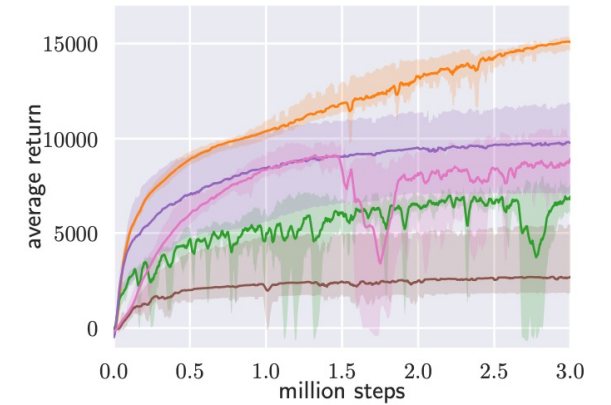


(a) Hopper-v1

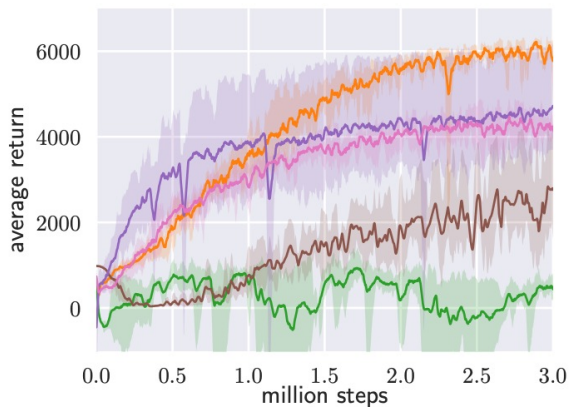


(b) Walker2d-v1

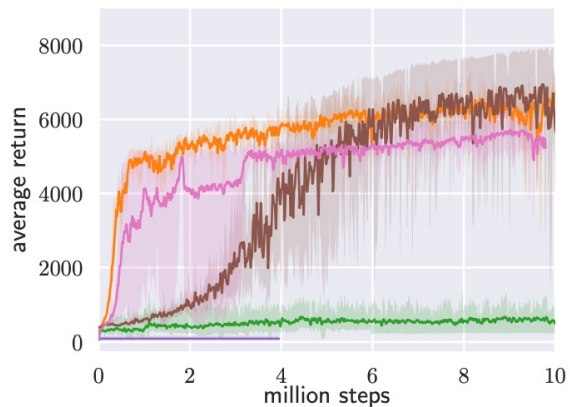
From Haarnoja, Zhou et al. (2018)



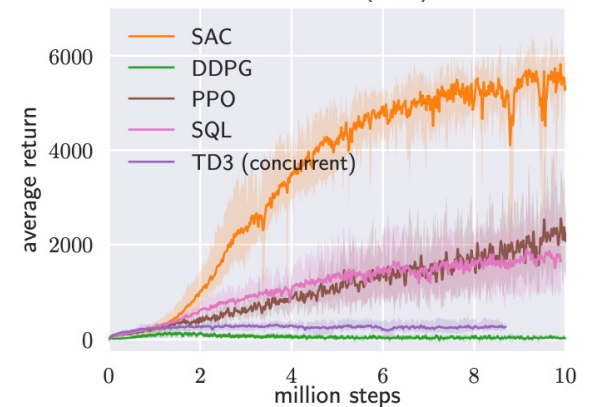
(c) HalfCheetah-v1



(d) Ant-v1



(e) Humanoid-v1



(f) Humanoid (rllab)

Robustness to Environment Changes

SAC on Minotaur - Testing

Using Soft Actor Critic (SAC), Minotaur learns to walk quickly and to generalize to environments with challenges that it was not trained to deal with!

