

# Soft Q-learning with Mutual Information Regularization

Accepted in ICLR 2019

Kyeong Ryeol, Go  
M.S. Candidate of OSI Lab

# Motivation

- Entropy-regularization RL

$$\max_{\pi} E \left[ \sum_t \gamma^t \left( r(s_t, a_t) - \frac{1}{\beta} \log \pi(a_t | s_t) \right) \right]$$

- KL-regularized RL

$$\max_{\pi} E \left[ \sum_t \gamma^t \left( r(s_t, a_t) - \frac{1}{\beta} \log \pi(a_t | s_t) / \rho(a_t) \right) \right]$$

- \* When some actions are simply non-useful or not frequently used
- \* When actions have significantly different importance depending on the task

# Paper list

- When some actions are simply non-useful or not frequently used
  - Soft q-learning with mutual information regularization (ICLR 2019)
- When actions have significantly different importance
  - Information asymmetry in KL-regularized RL (ICLR 2019)
  - Exploiting hierarchy for learning and transfer in KL-regularized RL (Arxiv)

*One-step decision-making...*

$$\begin{aligned} & \max_{\pi, \rho} E \left[ r(s, a) - \frac{1}{\beta} \log \pi(a|s) / \rho(a) \right] \\ &= \max_{\pi} \sum_{s,a} p(s) \pi(a|s) r(s, a) - \frac{1}{\beta} \min_{\rho} \sum_s p(s) KL(\pi(\cdot | s) \parallel \rho(\cdot)) \\ &= \max_{\pi} \sum_{s,a} p(s) \pi(a|s) r(s, a) - \frac{1}{\beta} \mathbf{I}(\mathcal{S}, \mathcal{A}) \quad \left( \because \rho^*(a) = \sum_s p(s) \pi(a|s) \right) \\ & \qquad \qquad \text{mutual-information regularization} \end{aligned}$$

Also hold in *Multi-step decision-making(=RL)* when  $\gamma \rightarrow 1$  (See Appendix)

1. What is an optimal policy for a fixed action prior?
2. What is an optimal action prior for a fixed policy?

- Definitions

- Transition probability

- $P_{\pi}^t(s'|s) = \sum_a P(s'|a, s)\pi(a|s) \quad s.t. \quad \mathbf{P}_{\pi} \in \mathbb{R}^{|\mathcal{S}|} \times \mathbb{R}^{|\mathcal{S}|}$

- Stationary distribution over states (assumed to exist)

- $\boldsymbol{\mu}_{\pi}^T := \lim_{t \rightarrow \infty} \mathbf{v}_0^T \mathbf{P}_{\pi}^t \quad s.t. \quad \mu_{\pi}(s') = \sum_s P_{\pi}(s'|s)\mu_{\pi}(s), \quad \boldsymbol{\mu}_{\pi}^T = \boldsymbol{\mu}_{\pi}^T \mathbf{P}_{\pi}$

- Stationary distribution over actions

- $\rho_{\pi}(a) := \sum_s \mu_{\pi}(s)\pi(a|s)$

1. What is an optimal policy for a fixed action prior?
2. What is an optimal action prior for a fixed policy?

$$\max_{\pi} E \left[ \sum_t r(s_t, a_t) - \frac{1}{\beta} \log \pi(a_t | s_t) / \rho(a_t) \right]$$

- $V_{\pi, \rho}(s) := E \left[ \sum_t \gamma^t \left( r(s_t, a_t) - \frac{1}{\beta} \log \pi(a_t | s_t) / \rho(a_t) \right) \mid s_0 = s \right]$
- $Q_{\pi, \rho}(s, a) := r(s, a) + \gamma E_{s'} [V_{\pi, \rho}(s')]$

By standard variational calculus,

$$\begin{aligned} \pi^*(a|s) &\propto \rho(a) \exp \left( \beta Q_{\pi^*, \rho}(s, a) \right) \\ &\rightarrow \pi^*(a|s) \propto \exp \left( \beta Q_{\pi^*, \rho}(s, a) \right) \text{ if } \rho(a) = 1/|\mathcal{A}| \\ &\rightarrow \pi^*(a|s) = 1 \text{ if } a = \arg \max_a Q^*(s, a) \text{ if } \beta \rightarrow \infty \end{aligned}$$

1. What is an optimal policy for a fixed action prior?
2. What is an optimal action prior for a fixed policy?

$$\begin{aligned}
 & \arg \max_{\rho} E \left[ \sum_t \gamma^t \left( r(s_t, a_t) - \frac{1}{\beta} \log \pi(a_t | s_t) / \rho(a_t) \right) \right] \\
 &= \arg \max_{\rho} \sum_t \sum_s \gamma^t \mathbf{v}_t(s) \sum_a \pi(a | s) \left( r(s, a) - \frac{1}{\beta} \log \pi(a | s) / \rho(a) \right) \\
 &= \arg \max_{\rho} - \frac{1}{\beta} \sum_t \sum_s \gamma^t \mathbf{v}_t(s) KL(\pi(\cdot | s) \parallel \rho(\cdot))
 \end{aligned}$$

- $\mathbf{v}_t(s) := \sum_{s_0, a_0, \dots, s_{t-1}, a_{t-1}} p(s_0) \left( \prod_{t'=0}^{t-2} \pi(a_{t'} | s_{t'}) P(s_{t'+1} | s_{t'}, a_{t'}) \right) \pi(a_{t-1} | s_{t-1}) P(s | s_{t-1}, a_{t-1})$

$$\begin{array}{c}
 \rho^*(a) \propto \sum_s \sum_T \gamma^t \mathbf{v}_t(s) \pi(a | s) \\
 \downarrow \qquad \qquad \qquad \searrow \\
 \text{stationary distribution over actions} \qquad \mu_{\pi}(s) : \text{stationary distribution over states}
 \end{array}$$

# MIRL : Mutual Information RL

## 1. Tabular setting

- Q-functions update

$$Q(s, a) \leftarrow (1 - \alpha_Q)Q(s, a) + \alpha_Q \left( T_{soft}^\rho Q \right) (s, a, s')$$
$$\text{where } \left( T_{soft}^\rho Q \right) (s, a, s') := r(s, a) + \gamma \frac{1}{\beta} \log \sum_{a'} \rho(a') \exp(\beta Q(s', a'))$$

Why?

- $V_{\pi, \rho}(s) = E_a[r(s, a) - \frac{1}{\beta} \log \pi(a|s)/\rho(a) + \gamma E_{s'}[V_{\pi, \rho}(s')]]$
- $Q_{\pi, \rho}(s, a) = r(s, a) + \gamma E_{s', a'} \left[ Q_{\pi, \rho}(s', a') - \frac{1}{\beta} \log \pi(a'|s')/\rho(a') \right]$

Substitute  $\pi(a|s)$  to  $\pi^{upd}(a|s) \propto \rho(a) \exp(\beta Q_{\pi, \rho}(s, a))$



# MIRL : Mutual Information RL

## 1. Tabular setting

- Prior update

$$\rho_{i+1}(a) = (1 - \alpha_\rho)\rho_i(a) + \alpha_\rho\pi_i(a|s_i)$$

where  $s_i \sim \nu_i(s)$  and  $\pi_i(a|s_i) \propto \rho_i(a) \exp(\beta Q_i(s_i, a))$

↓ Converge to...

$$\rho_{\pi_i}(a) = \sum_s \nu_i(s)\pi_i(a|s) \left( = \sum_s \mu_{\pi_i}(s)\pi_i(a|s) \text{ with } \gamma = 1 \right)$$

Common practice  
in actor-critic

- $\beta$  update

$$\beta_i = c \cdot i$$

# MIRL : Mutual Information RL

## 2. High-dimensional state space

- Q-functions update

$$L(\theta, \rho) := E_{s,a,r,s' \sim \mathcal{M}} \left[ \left( \left( T_{soft}^{\rho} Q_{\bar{\theta}} \right) (s, a, s') - Q_{\theta}(s, a) \right)^2 \right]$$

- Prior update

$$\rho_{i+1}(a) = (1 - \alpha_{\rho})\rho_i(a) + \alpha_{\rho}\pi_i(a|s_i)$$

- $\beta$  update

$$\beta_{i+1} = (1 - \alpha_{\beta})\beta_i + \alpha_{\beta} \frac{1}{L(\theta_i, \rho_{i+1})}$$

---

**Algorithm 1** MIRL

---

```
1: Input: the learning rates  $\alpha_\rho$ ,  $\alpha_Q$  and  $\alpha_\beta$ , a Q-network  $Q_\theta(s, a)$ , a target network  $Q_{\bar{\theta}}(s, a)$ , a  
   behavioural policy  $\pi_b$ , an initial prior  $\rho_0$  and parameters  $\theta_0$  at  $t = 0$ .  
2: for  $i = 1$  to  $N$  iterations do  
3:   Get environment state  $s_i$  and apply action  $a_i \sim \pi_b(\cdot|s_i)$   
4:   Get  $r_i, s_{i+1}$  and store  $(s_i, a_i, r_i, s_{i+1})$  in replay memory  $\mathcal{M}$   
5:   Update prior  $\rho_{i+1}(\cdot) = \rho_i(\cdot)(1 - \alpha_\rho) + \alpha_\rho \pi_i(\cdot|s_i)$   
6:   if  $i \bmod \text{update frequency} == 0$  then  
7:     Update Q-function  $\theta_{i+1} = \theta_i - \alpha_Q \nabla_\theta L(\theta, \rho_{i+1})|_{\theta_i}$  according to Equation (13)  
8:     Update parameter  $\beta_{i+1} = (1 - \alpha_\beta)\beta_i + \alpha_\beta \left( \frac{1}{L(\theta_i, \rho_{i+1})} \right)$   
9:   end if  
10: end for
```

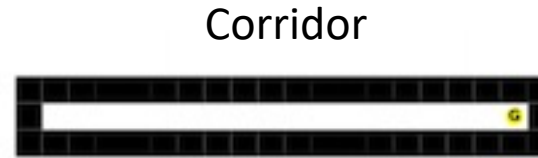
---

- Behavior policy  $\pi_b(a|s_i)$

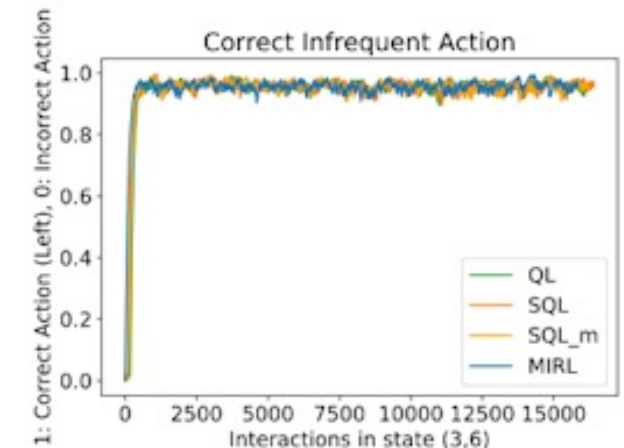
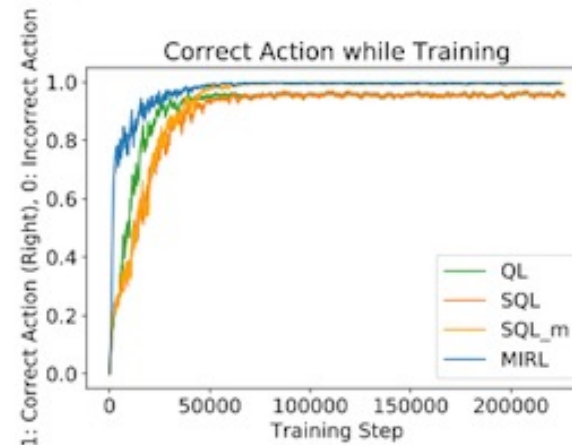
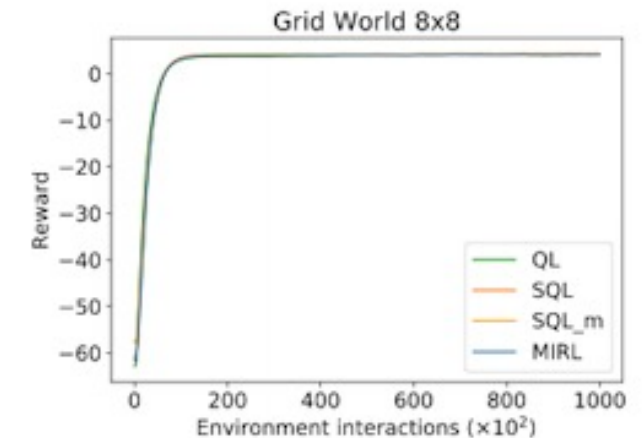
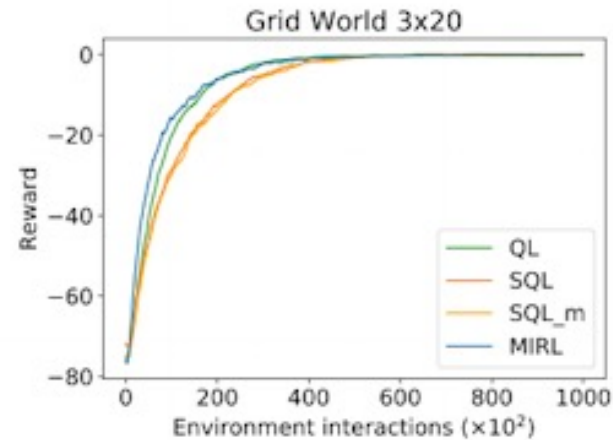
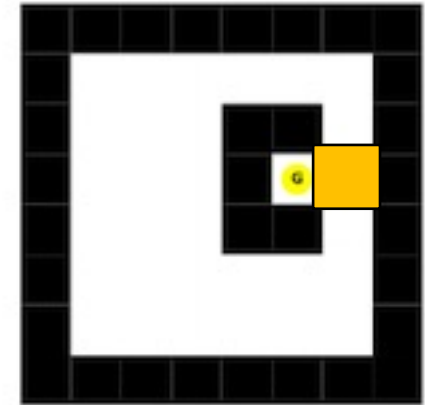
$$a_i = \begin{cases} \arg \max_a \pi_i(a|s_i) & \text{if } u \geq \epsilon & \text{Exploitation} \\ a \sim \rho_i(a) & \text{if } u \leq \epsilon & \text{Exploration} \end{cases}$$

# Grid-world

- Settings
  - Reaching a goal : reward 9
  - Else : reward -1
  - Restarted in a random location
- Baseline
  - > Q-learning (QL)
  - > Soft Q-learning (SQL)
  - > Soft Q-learning + behavior policy (SQL\_m)
  - > MIRL

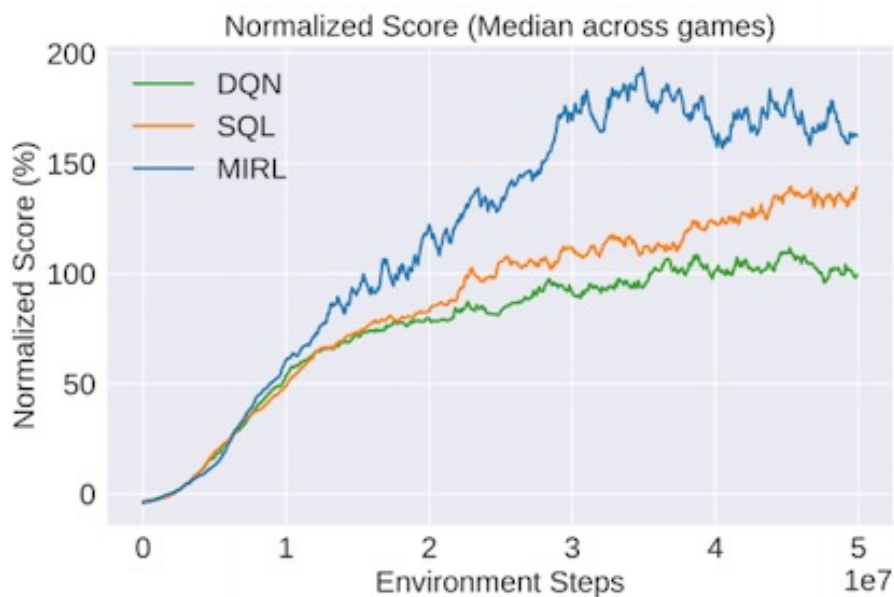


Square world



# ATARI

- Baselines
  - > Deep Q-Network (DQN)
  - > Soft Q-learning (SQL)
  - > MIRL



$$Z_{normalized} = \frac{Z - Z_{random}}{Z_{human} - Z_{random}} \times 100\%$$

Game	DQN (%)	SQL (%)	MIRL (%)
Alien	<b>101.58</b>	51.02	40.23
Assault	250.61	283.62	<b>357.40</b>
Asterix	166.32	242.73	<b>330.19</b>
Asteroids	<b>9.74</b>	8.57	7.80
BankHeist	97.12	94.62	<b>166.26</b>
BeamRider	99.16	113.64	<b>117.21</b>
Boxing	2178.57	2283.33	<b>2338.89</b>
ChopperCommand	<b>72.71</b>	26.37	65.03
DemonAttack	350.95	451.78	<b>469.30</b>
Gopher	474.18	<b>538.87</b>	429.44
Kangaroo	351.48	393.16	<b>405.9</b>
Krull	843.16	886.68	<b>1036.04</b>
KungFuMaster	122.14	<b>142.04</b>	121.41
Riverraid	77.21	<b>109.37</b>	76.02
RoadRunner	548.90	613.62	<b>695.88</b>
Seaquest	21.95	36.00	<b>64.86</b>
SpaceInvaders	166.62	<b>200.38</b>	164.79
StarGunner	653.44	<b>681.12</b>	574.89
UpNDown	183.19	230.82	<b>394.21</b>
Mean	356.26	388.83	<b>413.46</b>

Table 1: Mean Normalized score in 19 Atari games for DQN, SQL and our approach MIRL.



# Ablation study

- Baselines
  - > Soft Q-learning (SQL)
  - > Soft Q-learning + behavior policy (SQL\_m)
  - > MIRL

