

Standard MDP

State :	S
Actions :	A
Reward :	$r(s,a)$
Transition :	$p(s' s,a)$
Discount :	γ
optimal policy:	$\pi^* = \arg\max_{\pi} \sum_t \gamma^t E_{s,a \sim \pi} [r(s,a)]$

Maximum Entropy MDP

	S
	A
	$r(s,a) + \alpha H(\pi(s))$
	$p(s' s,a)$
	γ
	$\pi_{\text{soft}}^* = \arg\max_{\pi} \sum_t \gamma^t E_{s,a \sim \pi} [r(s,a) + \alpha H(\pi(s))]$

Definition:

$$\text{optimal policy: } \pi_{\text{soft}}^* = \arg\max_{\pi} \sum_t \gamma^t E_{s,a \sim \pi} [r(s,a) + \alpha H(\pi(s))]$$

$$\text{value function: } V_{\text{soft}}^*(s) = \sum_t \gamma^t E_{s,a \sim \pi} [r(s,a) + \alpha H(\pi(s))]$$

$$\text{Action-value function: } Q_{\text{soft}}^*(s,a) = r(s,a) + \sum_t \gamma^t E_{s,a \sim \pi} [r(s,a) + \alpha H(\pi(s))]$$

Theorem

$$Q_{\text{soft}}^*(s,a) = r(s,a) + \gamma E_s [V_{\text{soft}}^*(s)]$$

$$V_{\text{soft}}^*(s) = \alpha H(\pi(s)) + \sum_a \pi(a|s) Q_{\text{soft}}^*(s,a) \quad \text{Easy to see}$$

Theorem: Energy Form

Given a policy π with value function $V_{\text{soft}}^*, V_{\text{soft}}$

the greedy policy wrt π is $\tilde{\pi}(s) = \arg\max_{\pi(s)} \sum_a \pi(a|s) Q_{\text{soft}}^*(s,a) + \alpha H(\tilde{\pi}(s))$

$$Q_{\text{soft}}^*(s, \pi(s)) = \alpha \log \sum_a \exp(Q_{\text{soft}}^*(s,a)/\alpha)$$

$$\tilde{\pi}(a|s) = \frac{\exp(Q_{\text{soft}}^*(s,a)/\alpha)}{\sum_a \exp(Q_{\text{soft}}^*(s,a)/\alpha)} = \exp\left[\frac{1}{\alpha} [Q_{\text{soft}}^*(s,a) - Q_{\text{soft}}^*(s,\pi(s))]\right]$$

$$\tilde{\pi}(s) = \arg\max_{\pi(s)} \sum_a \pi(a|s) Q_{\text{soft}}^*(s,a) + \alpha H(\tilde{\pi}(s))$$

$$= \arg\max_{\pi(s)} \sum_a \pi(a|s) [Q_{\text{soft}}^*(s,a) - \alpha \log \tilde{\pi}(a|s)]$$

$$\begin{aligned} \min_{\pi} & - \sum_i \chi_i (a_i - \alpha \log \tilde{\pi}(a_i)) \\ \text{s.t.} & \sum_i \chi_i = 1 \end{aligned} \quad \Rightarrow \text{Convex} \quad \begin{array}{l} (\text{optimal value is negated}) \\ \text{when flip max to min} \end{array}$$

$$L(x_i) = -\sum_j x_j (a_{ij} - \alpha(\gamma x_i)) + V\left(\sum_j x_j - 1\right)$$

$$\frac{\partial L}{\partial x_i} = -a_{ii} + \alpha(1 + \gamma x_i) + V := 0$$

$$x_i^*(v) = \exp\left(\frac{a_{ii} - v}{\alpha} - 1\right)$$

$$g(v) = -\sum_i \exp\left(\frac{a_{ii} - v}{\alpha} - 1\right)(v + \alpha) + V\left(\sum_i \exp\left(\frac{a_{ii} - v}{\alpha} - 1\right) - 1\right)$$

$$= -\alpha \sum_i \exp\left(\frac{a_{ii} - v}{\alpha} - 1\right) - V$$

$$\frac{\partial g}{\partial v} = \sum_i \exp\left(\frac{a_{ii} - v}{\alpha} - 1\right) - 1 = \exp\left(\frac{v}{\alpha}\right) \sum_i \exp\left(\frac{a_{ii}}{\alpha} - 1\right) - 1 := 0$$

$$v^* = \alpha \log \sum_i \exp\left(\frac{a_{ii}}{\alpha} - 1\right)$$

$$x_i^* = \frac{\exp\left(\frac{a_{ii}}{\alpha} - 1\right)}{\exp\left(\frac{v^*}{\alpha}\right)} = \frac{\exp\left(\frac{a_{ii}}{\alpha} - 1\right)}{\sum_j \exp\left(\frac{a_{ij}}{\alpha} - 1\right)} = \frac{\exp(a_{ii}/\alpha)}{\sum_j \exp(a_{ij}/\alpha)}$$

$$p^* = -\sum_i \frac{\exp(a_{ii}/\alpha)}{\sum_j \exp(a_{ij}/\alpha)} \left(a_{ii} - \alpha \log \frac{\exp(a_{ii}/\alpha)}{\sum_j \exp(a_{ij}/\alpha)} \right)$$

$$= -\frac{1}{\alpha} \sum_i \exp(a_{ii}/\alpha) (\alpha \log z)$$

$$= -\alpha \log z$$

$$= -\alpha \log \sum_i \exp(a_{ii}/\alpha)$$

$$p^* = -\alpha \log \sum_i \exp(a_{ii}/\alpha)$$

$$x_i^* = \exp\left(\frac{a_{ii}}{\alpha} + \frac{p^*}{\alpha}\right)$$

Theorem: Monotonic Improvement

given a policy π with value function Q_{soft}^{π}

let one-step greedy policy $\tilde{\pi}(a|s) = \frac{\exp(Q_{\text{soft}}^{\pi}(s,a)/\alpha)}{\sum_a \exp(Q_{\text{soft}}^{\pi}(s,a)/\alpha)}$

then $Q_{\text{soft}}^{\tilde{\pi}}(s,a) > Q_{\text{soft}}^{\pi}(s,a) \forall (s,a)$

thus the optimal policy must have form $\pi^*(a|s) = \frac{\exp(Q_{\text{soft}}^*(s,a)/\alpha)}{\sum_a \exp(Q_{\text{soft}}^*(s,a)/\alpha)}$

$$Q_{\text{soft}}^{\pi}(s,a) = E_{s'} [r_o + \gamma (Q_{\text{soft}}^{\pi}(s',a') + E_{a' \sim \pi} [Q_{\text{soft}}^{\pi}(s',a')])]$$

$$\leq E_{s'} [r_o + \gamma (Q_{\text{soft}}^{\tilde{\pi}}(s') + E_{a' \sim \tilde{\pi}} [Q_{\text{soft}}^{\tilde{\pi}}(s',a')])]$$

$\leq \dots$

$$\leq Q_{\text{soft}}^{\tilde{\pi}}(s,a)$$

Theorem: Bellman Optimality

By the monotonic improvement theorem, $\pi^*(a|s) = \frac{\exp(\alpha_{sa}^*(s,a))}{\sum_a \exp(\alpha_{sa}^*(s,a))}$

$$\text{then } \pi^* = \arg\max_{\pi} \sum_a \pi(a|s) \alpha_{sa}^*(s,a) + \gamma \mathbb{E}[\pi(s)]$$

$$V_{\pi^*}(s) = \sum_a \pi^*(a|s) \alpha_{sa}^*(s,a) + \gamma \mathbb{E}[\pi^*(s)]$$

$$= \max_{\pi} \sum_a \pi(a|s) \alpha_{sa}^*(s,a) + \gamma \mathbb{E}[\pi(s)]$$

$$= \alpha \log \sum_a \exp(\alpha_{sa}^*(s,a)/\alpha)$$

$$\alpha_{sa}^*(s,a) = r(s,a) + \gamma \mathbb{E}_s [V_{\pi^*}(s)]$$

$$= r(s,a) + \gamma \mathbb{E}_s [\alpha \log \sum_a \exp(\alpha_{sa}^*(s,a)/\alpha)]$$

At each state, the optimal policy π^*
is greedy wrt α_{sa}^*

Theorem: Bellman Backup and uniqueness of π^*

define the soft value iteration as

$$(TQ)(s,a) = r(s,a) + \gamma \mathbb{E}_s [\alpha \log \int_{a'} \exp(\alpha_{sa}^*(s,a')/\alpha) da']$$

Then T is a γ -contraction

Thus there's a unique α_{sa}^* satisfying

$$\alpha^*(s,a) = r(s,a) + \gamma \mathbb{E}_s [\alpha \log \int_{a'} \exp(\alpha_{sa}^*(s,a')/\alpha) da']$$

and the optimal policy is $\pi^*(a|s) = \frac{\exp(\alpha_{sa}^*(s,a)/\alpha)}{\int_a \exp(\alpha_{sa}^*(s,a)/\alpha) da}$

$$\text{let } \| \alpha_1 - \alpha_2 \|_\infty = \max_{s,a} | \alpha_1(s,a) - \alpha_2(s,a) | = \epsilon$$

$$\begin{aligned} \alpha \log \int_a \exp \frac{\alpha_1(s,a)}{\alpha} da &\leq \alpha \log \int_a \exp \frac{\alpha_2(s,a) + \epsilon}{\alpha} da \\ &= \alpha \log \int_a \exp \frac{\epsilon}{\alpha} \exp \frac{\alpha_2(s,a)}{\alpha} da \\ &= \epsilon + \alpha \log \int_a \exp \frac{\alpha_2(s,a)}{\alpha} da \quad \forall s \end{aligned}$$

$$\alpha \log \int_a \exp \frac{\alpha_1(s,a)}{\alpha} da \geq -\epsilon + \alpha \log \int_a \exp \frac{\alpha_2(s,a)}{\alpha} da \quad \forall s$$

$$-\epsilon \leq \alpha \log \int_A \exp \frac{\alpha_i(s,a)}{\alpha} da - \alpha \log \int_A \exp \frac{\alpha_i(s,a)}{\alpha} da \leq \epsilon \quad \forall s$$

$$\begin{aligned}\|T\alpha_1 - T\alpha_0\|_\infty &= \max_{s,a} |E_s \left[\alpha \log \int_A \exp \frac{\alpha_i(s,a)}{\alpha} da \right] - E_{s'} \left[\alpha \log \int_A \exp \frac{\alpha_i(s,a)}{\alpha} da \right]| \\ &= \max_{s,a} |f \left[\alpha \log \int_A \exp \frac{\alpha_i(s,a)}{\alpha} da \right] - \alpha \log \int_A \exp \frac{\alpha_i(s,a)}{\alpha} da| \\ &\leq f\epsilon \\ &= \epsilon \| \alpha_1 - \alpha_0 \|_\infty\end{aligned}$$

• Soft Q-Iteration

$$\begin{aligned}\alpha_{soft}(s,a) &\leftarrow r(s,a) + \gamma \cdot \alpha \log \int_A \exp \frac{\alpha(s,a')}{\alpha} da' \xrightarrow{\text{for discrete action space}} \alpha \log \sum_a \exp \frac{\alpha(s,a)}{\alpha} da \\ &= r(s,a) + \gamma \alpha \log E_{g(a')} \left[\frac{\alpha(s,a')}{g(a')} \right]\end{aligned}$$

Initialize α^θ and target $\bar{\alpha} = \alpha^\theta$

initialize replay buffer D

while true

Sample $(s,a,r,s') \sim D$ $\xrightarrow{\text{can be any distribution, for example uniform}}$

Sample $a'_i \sim g(a')$ $\xrightarrow{\text{but scales poorly to high dimension}}$

Compute $t = r + \gamma \log \left[\frac{1}{m} \sum_{a'_i} \frac{\exp(\bar{\alpha}(s',a'_i)/\alpha)}{g(a'_i)} \right]$ one choice is $g(a'|s) = \frac{\exp(\alpha(s,a)/\alpha)}{\int_a \exp(\alpha(s,a)/\alpha) da}$

$\min_\theta \| \alpha^\theta(s,a) - t \|^2$

for each k iteration: $\bar{\alpha} \leftarrow \alpha$

$\xrightarrow{\text{but this will require SGD for continuous action space}}$

• Soft Actor Critic (Standard)

Compute α , π_θ , V_θ

$$\text{Policy } \pi_\theta(a|s) = \frac{\exp(\alpha_\theta(s,a)/\kappa)}{\int_a \exp(\alpha_\theta(s,a)/\kappa) da}$$

$$\min_{\phi} D_{KL}\left(\pi_\phi(\cdot|s) \parallel \frac{\exp(\alpha_\phi(s,a)/\kappa)}{\int_a \exp(\alpha_\phi(s,a)/\kappa) da}\right) = \min_{\phi} \int_a \pi_\phi(a|s) \cdot \log \frac{\pi_\phi(a|s)}{\exp(\alpha_\phi(s,a)/\kappa) / \int_a \exp(\alpha_\phi(s,a)/\kappa) da} da$$

$$\min_{\phi} E_{a \sim \pi_\phi(s)} [\log \pi_\phi(a|s) - \alpha_\phi(s,a)/\kappa] = \min_{\phi} E_{a \sim \pi_\phi(s)} [\log \pi_\phi(a|s) - \alpha_\phi(s,a)/\kappa]$$

$$\text{let } \pi_\phi(s) = \mathbb{U}_\phi(s) + \alpha_\phi(s) \cdot \mathcal{G} \quad \mathcal{G} \sim N(0, I)$$

$$\min_{\phi} E_{s \sim \mathcal{D}} \left[\alpha \log \pi_\phi(\mathbb{U}_\phi(s) + \alpha_\phi(s) \cdot \mathcal{G}) - \alpha_\phi(s, \mathbb{U}_\phi(s) + \alpha_\phi(s) \cdot \mathcal{G}) \right]$$

Initialize V_ψ , $V_\theta = V_\psi$, α_θ , α_ϕ , π_θ

while true {

Sample $(s, a, r, s') \sim \mathcal{D}$

learn α to approximate $\alpha_{soft}^*(s, a) = r(s, a) + \gamma E_s [V_{soft}^*(s')]$

compute target $y(s) = r(s, a) + \gamma V_\psi(s')$

$$\min_{\alpha_1} \|\alpha_{\theta_1}(s, a) - y(s)\|_2 \quad \min_{\alpha_2} \|\alpha_{\phi_2}(s, a) - y(s)\|_2$$

Sample $\mathcal{G} \sim N(0, I) \quad \mathcal{A} = \mathbb{U}_\phi(s) + \alpha_\phi(s) \cdot \mathcal{G}$

$$\begin{aligned} \text{learn } V \text{ to approximate } V_{soft}^*(s) &= \alpha \mathbb{U}_\phi(\mathbb{U}_\phi(s)) + \int_a \pi_\phi(a|s) \cdot \alpha_{soft}^*(s, a) da \\ &= E_{a \sim \pi_\phi(s)} [\alpha_{soft}^*(s, a) - \alpha \log \pi_\phi(a|s)] \end{aligned}$$

$$\text{compute target } y(s) = \min (\alpha_\theta(s, a), \alpha_\phi(s, a)) - \alpha \log \pi_\phi(a|s)$$

$$\min_{\psi} \|V_\psi(s) - y(s)\|_2$$

$$\text{learn } \pi \text{ to approximate } \pi^*(a|s) = \frac{\exp(\alpha^*(s, a)/\kappa)}{\int_a \exp(\alpha^*(s, a)/\kappa) da}$$

$$\min_{\phi} \alpha \log \pi_\phi(a|s) - \min (\alpha_\theta(s, a), \alpha_\phi(s, a))$$

$$\bar{\psi} = \rho \bar{\psi} + (1-\rho) \psi$$

}

Soft Actor-Critic (Modem)

Initialize $\theta_1 = \theta_{\text{ini}}$ $\theta_2 = \theta_{\text{ini}}$ π_θ

while true {

sample $(s, a, r, s') \sim D$

$$\begin{aligned} \text{learn } \theta \text{ to approximate } Q_{\text{soft}}^*(s, a) &= r(s, a) + \gamma E_s [V_{\pi_\theta}^*(s)] \\ &= r(s, a) + \gamma E_{s'} \left[\int_a \pi^*(a'|s) Q_{\text{soft}}^*(s, a') + \alpha \log \pi^*(a'|s) \right] \\ &= r(s, a) + \gamma E_{s'} E_{a' \sim \pi(s)} \left[Q_{\text{soft}}^*(s, a') - \alpha \log \pi^*(a'|s) \right] \end{aligned}$$

$$\text{Sample } a' \sim \pi_\theta(a'|s) \sim u_\theta(s) + b_\theta(s) + \epsilon$$

$$\text{Compute target } y(s) = r(s, a) + \gamma \left[\min_{\theta_1} (\theta_1(s, a') - \alpha \log \pi_\theta^*(a'|s)) \right]$$

$$\min_{\theta_1} \|Q_{\theta_1}(s, a) - y(s)\|_2^2 \quad \min_{\theta_2} \|Q_{\theta_2}(s, a) - y(s)\|_2^2$$

$$\text{learn } \pi_\theta \text{ to approximate } \pi^*(a|s) = \frac{\exp(Q^*(s, a) / \alpha)}{\int_a \exp(Q^*(s, a) / \alpha) da}$$

Sample $a \sim \pi_\theta(a|s)$

$$\min_{\phi} \log \pi_\theta(a|s) - \min_{\theta_1} (\theta_1(s, a), \theta_2(s, a))$$

$$\bar{\theta}_1 = p\bar{\theta}_1 + (1-p)\theta_1$$

$$\bar{\theta}_2 = p\bar{\theta}_2 + (1-p)\theta_2$$

}

• Automating Entropy Adjustment

An Intuition (Based on dual gradient method explanation in EE364B
the dual gradient is basically adaptive penalty for violation of constraint)

Constraint: $E_{\alpha \sim \Delta} [\log \pi_\alpha(s)] \leq -\bar{f}_t : \alpha$

update $\alpha \leftarrow \underbrace{E_{\alpha \sim \Delta} [\log \pi_\alpha(s)] + \bar{f}_t}_{\rightarrow \bar{f}_t = -\dim A}$

if this is positive, the constraint is violated
then penalty is increased