

March 23rd (T) YY, ZHL.

1. Demo {
 a Photo/Image
 b Video

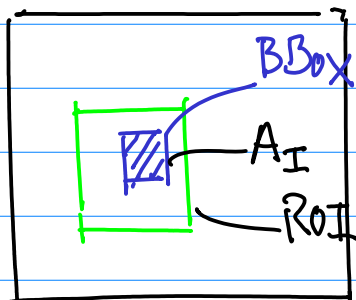
~ ravity — main Program

~ Box.py —
 Bounding

Testing Program: Rectangle_intersection.py.

Todo: Send readme.txt,

Threshold $T = 0.25$



Area (Intersection)
 Blue

A_I v.s. T

$A_I \geq T \rightarrow$ Alarm
 $A_I < T$ No Action

To Learn Policy

$$\pi(\vec{a}_t | \vec{s}_t) \quad \dots (2)$$

to maximize Eqn (1).

Augment Eqn (1) with
 Entropy to maximize
 Entropy at each visited
 State.

Define 2 Thresholds {
 Upper Threshold T_u
 Lower Threshold T_L

A_I
 $T_u \approx 50\%$
 A_I Interested
 $T_L \approx 10\%$
 Discarded
 A_I is not to be considered

So
 $\propto \frac{1}{\pi(a_t | s_t)} \quad \dots (3)$
 Entropy Function
 Independent Variable π
 Policy
 Weight
 ("Temperature" parameter
 determines the relative
 importance of Entropy term)

PART II

CTI

March 24 (Wed)

DRL-SAC

Soft Actor Critic Algorithm & Apps.

UC Berkeley & Google PP. 4

Standard RL Objective to maximize
 Reward

$$\max \left\{ \sum_{\tau} E_{(s_t, a_t) \sim p_{\pi}} [r(\vec{s}_t, \vec{a}_t)] \right\} \quad \dots (1)$$

Discard
 A_I if it is Bigger than T_u

Lecture Note on SAC

CTI
PART II

14L March 25

2

Hence Reward in Eqn(1) becomes

$$\max \left\{ \sum_{\tau} E_{(s_t, a_t) \sim p_{\pi}} [r(\vec{s}_t, \vec{a}_t)] + \alpha \gamma I(\pi(a_t | s_t)) \right\}$$

Soft Policy Iteration

Q-value (Reward)

$$Q: S \times A \rightarrow \mathbb{R}$$

... (1)

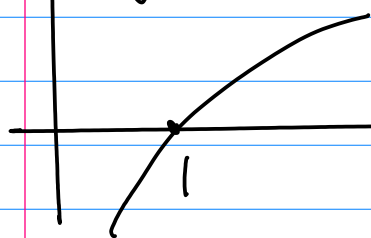
$$T^{\pi} Q(\vec{s}_t, \vec{a}_t) \triangleq r(\vec{s}_t, \vec{a}_t) + \gamma E_{\vec{s}_{t+1} \sim p} [V(\vec{s}_{t+1})] \dots (1^*)$$

Bellman
operator
for each action
leads to Reward

where

$$V(\vec{s}_{t+1}) = E_{\vec{a}_t \sim \pi} [Q(\vec{s}_t, \vec{a}_t) - \alpha \log \pi(\vec{a}_t | \vec{s}_t)] \dots (2)$$

$\log \pi(\cdot)$



Note: Based on (2)

$\pi(\cdot) \propto \log \pi(\vec{a}_t | \vec{s}_t)$, so $\pi(\vec{a}_t | \vec{s}_t)$ is

Reward Based Function

Policy introduced Q-Value
(Reward)

Lemma 1. Bellman Operator T "0" initial time $t=0$.

Given initial Condition $\vec{Q}^0: S \times A \rightarrow \mathbb{R}$ with $|A| < \infty$

Define Bellman Operator T as

$$Q^{k+1} = T^{\pi} Q^k \dots (3)$$

"All" total Number of finite
Actions

Part II

CTI

3

$$Q^{k+1} = T^{\pi} Q^k$$

Time Index

Q-Value

(Reward) Function Bellman Operator under Policy π

Note: 1° "Argmin" find minimum from a collection of elements, functions, etc.

Example Argmin $\{1, -1.2, -110, 212\}$
 $= -110$

For $k=0$, we have

$$Q^1 = T^{\pi} Q^0$$

$$k=1, Q^2 = T^{\pi} Q^1 = T^{\pi} (T^{\pi} Q^0) = (T^{\pi})^2 Q^0 \dots (4)$$

...

$$k=i, Q^{i+1} = (T^{\pi})^{i+1} Q^0 \dots (5)$$

Update Q-Value (Reward)
 Update Policy π

$$Z^0 \pi(\vec{a}_t | \vec{s}_t) \text{ Policy}$$

$\pi(\cdot | \vec{s}_t)$ All policies for each/every \vec{a}_t from \vec{a}_t .

$\pi'(\cdot | \vec{s}_t)$ one selected from the All policies (all action, e.g. "o" for \vec{a}_t)

$$\pi_{\text{new}} = \underset{\pi' \in \Pi}{\text{Argmin}} D_{KL} \left(\pi'(\cdot | \vec{s}_t) \parallel \frac{e^{\frac{1}{2} Q^{\pi_{old}}(\vec{s}_t | \cdot)}}{Z^{\pi_{old}}(\vec{s}_t)} \right) \dots (6)$$

Note: 1° $\pi \in \Pi$
 Policy from (belongs to) all collections of the Policies
 Collection of All Policies

2° D_{KL} : Kull Back - Leibler Divergence

3° π' One π realization from the collection of All Policies, A particular Policy

Generalize Eqn(b), we have

Lemma 2.

$$\pi_{\text{New}} = f(\pi_{\text{old}} \text{ from one of } | \text{ Conditions } | \text{ the all possible Policies } \dots (bb))$$

$$Q^{\pi^*}(\vec{s}_t, \vec{a}_t) \geq Q^{\pi}(\vec{s}_t, \vec{a}_t)$$

where π^* is from Eqn(b)

f to be minimization problem, so

$$\pi_{\text{New}} = \underset{\pi}{\operatorname{argmin}} D_{\text{KL}}(\pi_{\text{old}} \text{ from one of } | \text{ Condition } | \text{ the all policies } \dots (bc))$$

Policy $\pi \in \Pi$, under Eqn(b) converges to

Probability or likelihood

Note: Eqn(b) needs Better Explanation, as why it is formulated as

Note: Parameterized Policy as Gaussian Distribution

$$\frac{e^{-\frac{1}{2} Q^{\pi_{\text{old}}}(\vec{s}_t, \cdot)}}{Z^{\pi_{\text{old}}}(\vec{s}_t)} \dots (b)$$

D_{KL} minimization Problem.

Now, introduce parameterized Q function, and Policy π .

Note: In NN Softmax Activation Function

We map Neuron Output in $(-\infty, +\infty)$ to probabilistic distribution $[0, 1]$, e.g.

$$f: z \in (-\infty, +\infty) \rightarrow f(z) \in [0, 1]$$

where

$$f(z_i) = \frac{e^{z_i}}{\sum_{j=1}^N e^{z_j}} \dots (8)$$

$Q_{\theta}, \pi_{\phi}, \theta(\text{Theta}), \phi(\text{phi})$ are

ϕ : Policy parameters

$$J_Q(\theta) = E_{(\vec{s}_t, \vec{a}_t) \sim \theta} \left[\frac{1}{2} \left(Q_\theta(\vec{s}_t, \vec{a}_t) - \left(r(\vec{s}_t, \vec{a}_t) + \gamma E_{\vec{s}_{t+1} \sim p} [V_\theta(\vec{s}_{t+1})] \right) \right)^2 \right] \quad \dots (9)$$

Objective Function

$Q(\theta)$ Q-value (Reward) Function with Parameter θ (Theta)

Eqn (1*) Bellman Iteration

$$\nabla J_Q(\theta) = E_{(\vec{s}_t, \vec{a}_t) \sim \theta} \left\{ \left[Q_\theta(\vec{s}_t, \vec{a}_t) - \left(r(\vec{s}_t, \vec{a}_t) + \gamma E_{\vec{s}_{t+1} \sim p} [V_\theta(\vec{s}_{t+1})] \right) \right] \cdot \nabla Q_\theta(\vec{s}_t, \vec{a}_t) \right\} \quad \dots (10)$$

$$J_\pi(\phi) = E_{\vec{s}_t \sim D} \left[E_{\vec{a}_t \sim \pi_\phi} [\alpha \log(\pi_\phi(\vec{a}_t | \vec{s}_t)) - Q_\theta(\vec{s}_t, \vec{a}_t)] \right] \quad \dots (11)$$

$$J_\pi(\phi) = E_{\vec{s}_t \sim D} [E_{\vec{a}_t \sim \pi_\phi} [\cdot]]$$

For all states \vec{s}_t For all actions \vec{a}_t

Note:

$$\frac{d}{dx} \ln(x) = \frac{1}{x}, \quad \frac{d}{dx} \log_a(x) = \frac{1}{x} \frac{1}{\ln(a)} \quad \dots (12)$$

$$\pi(\vec{a}_t, \vec{s}_t) \quad \text{Eqn (1), PP. 2}$$

$$\vec{a}_t = f_\phi(\epsilon_t; \vec{s}_t) \quad \dots (13) \quad \text{Parameterized Policy?}$$

Action

Action Vector

Multi-dimensional

$$\vec{a}_t = (a_{t1}, a_{t2}, \dots, a_{tn})$$

$$J_\pi(\phi) = E_{\vec{s}_t \sim D, \epsilon_t \sim N} [\alpha \log \pi_\phi(f_\phi(\epsilon_t; \vec{s}_t) | \vec{s}_t) - Q_\theta(\vec{s}_t, f_\phi(\epsilon_t; \vec{s}_t))] \quad \dots (14)$$

Parameter for π

$\pi_\phi(f_\phi(\epsilon_t | \vec{s}_t) | \vec{s}_t)$ is from

$\pi_\phi(\vec{a}_{\phi t} | \vec{s}_t)$, or is from

$\pi_\phi(\vec{a}_t | \vec{s}_t)$ is from Eqn(2) PP, 2.

$$\hat{\nabla}_\phi J_\pi(\phi) = \nabla_\phi \alpha \log(\pi_\phi(\vec{a}_t | \vec{s}_t)) + \nabla_{\vec{a}_t} \alpha \log(\pi_\phi(\vec{a}_t | \vec{s}_t)) -$$

Note: Jacobian

Suppose $\vec{f}: \mathbb{R}^n \rightarrow \mathbb{R}^m$

$$\vec{x} \in \mathbb{R}^n, \vec{f}(\vec{x}) \in \mathbb{R}^m \quad \dots (16a)$$

$$\nabla f_1(\vec{x}) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} \\ \vdots \\ \frac{\partial f_1}{\partial x_n} \end{bmatrix} \quad \dots (16e)$$

Example: $f(x_1, x_2) = f(\vec{x}), \vec{x} \in \mathbb{R}^2$

$$\text{Then } \vec{f}(\vec{x}) = (f_1(\vec{x}), f_2(\vec{x})), \quad \vec{x} = (x_1, x_2) \quad \dots (16b)$$

$$\vec{f}(\vec{x}) \in \mathbb{R}^2. \quad \dots (16c)$$

$$\nabla^T f_1(\vec{x}) = \left(\frac{\partial f_1}{\partial x_1} \quad \dots \quad \frac{\partial f_1}{\partial x_n} \right) \quad \dots (16f)$$

Jacobian: Matrix of All its 1st order

Partial Derivatives.

Hence, Jacobian

$$\vec{J} \triangleq \left[\frac{\partial \vec{f}}{\partial x_1}, \dots, \frac{\partial \vec{f}}{\partial x_n} \right] = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \frac{\partial f_m}{\partial x_1} & & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \quad \dots (16d)$$

$$\vec{J} = \begin{bmatrix} \nabla^T f_1 \\ \vdots \\ \nabla^T f_m \end{bmatrix} \quad \dots (16g)$$

gradient of $f_1(\vec{x})$

$$\nabla f_1 =$$

Note: $\nabla f_1(\vec{x})$, if $f(\vec{x})$ changed to $f(\vec{x}; \theta)$

$\nabla_{\phi} f_1(\vec{x}; \phi)$, so we are taking Partial derivative wrt ϕ .

$\phi = (\phi_1, \phi_2, \dots, \phi_n)$, $\vec{\phi}$ is written

Simplified as ϕ .

Automating Entropy Adjustment

$$\max_{\pi_{0:T}} E_{\rho_{\pi}} \left[\sum_{t=0}^T r(\vec{s}_t, \vec{a}_t) \right] \quad \text{s.t.} \quad E_{(\vec{s}_t, \vec{a}_t) \sim \rho_{\pi}} [-\log(\pi_t(\vec{a}_t | \vec{s}_t))] \geq H, \forall t \dots (17)$$

... (18)

$\hookrightarrow r(\vec{s}_t, \vec{a}_t)$: reward function

$\hookrightarrow \sum_{t=0}^T r(\vec{s}_t, \vec{a}_t)$: Summation of reward function from $t=0$ to $t=T$;

$$\hookrightarrow \max \left[\sum_{t=0}^T r(\vec{s}_t, \vec{a}_t) \right]$$

maximize the reward function from $t \in [0, T]$ period.

$$\hookrightarrow \max E \left[\sum_{t=0}^T r(\vec{s}_t, \vec{a}_t) \right]$$

maximize expected reward function from $[0, T]$

Now Add Notation to detail it up

$$\max_{\pi_{0:T}} E \left[\sum_{t=0}^T r(\vec{s}_t, \vec{a}_t) \right]$$

\uparrow ρ_{π} under Policy π from time 0 to T. ρ_{π} under Policy, Policy π

So, we have Eqn (17).

$$-\log(\pi_t(\vec{a}_t | \vec{s}_t)) = \log \frac{1}{\pi_t(\vec{a}_t | \vec{s}_t)}$$

Dynamic Programming, Solving
for the policy backward through time

See Ref. PP. 7
"SAC" Paper

Rewrite Objective Function, Egn (17)

$$\max_{\pi_0} \{ E[r(\vec{s}_0, \vec{a}_0)] + \max_{\pi_1} \{ E[\dots] + \max_{\pi_T} E[r(\vec{s}_T, \vec{a}_T)] \} \} \dots (18)$$

$$\max_{\pi_0} \{ E[r(\vec{s}_0, \vec{a}_0)] + \max_{\pi_1} \{ E[r(\vec{s}_1, \vec{a}_1)] + \max_{\pi_2} \{ E[r(\vec{s}_2, \vec{a}_2)] \} \} \} \text{ for } t=2$$

$$\max_{\pi_2} \{ E[r(\vec{s}_2, \vec{a}_2)] + \max_{\pi_3} \{ E[r(\vec{s}_3, \vec{a}_3)] \} \} \text{ for } t=3$$

$$\max_{\pi_3} \{ E[r(\vec{s}_3, \vec{a}_3)] + \max_{\pi_4} \{ E[r(\vec{s}_4, \vec{a}_4)] \} \} \text{ for } t=4$$

...