

Lecture 04

March 27, 2024

Revision so far, TRPO, PPO

$$P(x_1), P(x_2) \dots P(x_n)$$

$$H = -\sum_i P \log P$$

Heads / tails

SAC

usual use of entropy as regulariser

why not include in the reward directly?

$$J(\pi) = \sum_i E[r(s_t, a_t) + \alpha H[\pi(\cdot | s_t)]]$$

"soft value"

$$v(s_t) = E[Q(s_t, a_t) - \log \pi(a_t | s_t)]$$

Then update becomes,

$$\pi_{\text{new}} = \arg \min_{\pi' \in \Pi} D_{\text{KL}} \left[\pi'(\cdot | s_t) \parallel \frac{\exp(Q^{\pi_{\text{old}}}(s_t, \cdot))}{Z^{\pi_{\text{old}}}(s_t)} \right]$$

→ Q^{π} gets updated arbitrarily

→ π is parameterised, so we find a π_{new} within the family of distributions that minimises KL

Expand as $\pi_{\text{new}} = \arg \min E \left[\log \pi' - \underbrace{Q^{\pi_{\text{old}}}}_{\text{in the sense of policy evaluation}} + \underbrace{\log Z^{\pi_{\text{old}}}}_{\text{constant}} \right]$

Parameterise by three networks,

$$J_v(\psi), J_q(\theta), J_{\pi}(\phi)$$

all interdependent

Summary, motivation for trust region methods

Value based methods \rightarrow Approximate $Q(s_t, a_t)$
 \rightarrow Easy to use
 \rightarrow Difficult to scale

Policy gradient \rightarrow Direct optimisation of $v_{\theta}(s)$
 \rightarrow gradients may be difficult to distinguish

Actor critics \rightarrow Policy gradient with baseline "Advantage"
 \rightarrow Baseline provided by value function
 $\rightarrow v_{\theta}(s)$
 $\rightarrow Q_{\theta}(s, a)$

DDPG \rightarrow Continuous control
 \rightarrow Removes need to compute expectation over a_{t+1} in $r_t + E_{a_{t+1}}[Q(s_{t+1}, a_{t+1})]$
 \rightarrow Requires addition of noise at the end for exploration

SAC \rightarrow Include entropy in the reward itself

Summary of challenges \rightarrow A lot of hyperparameters to tune
 \rightarrow Can we guarantee monotonic improvement for a certain step size?

TRPO

→ Defines cost rather than reward
 ⇒ positive advantage is bad

$$\underbrace{\eta(\tilde{\pi})}_{\text{cost function}} = \eta(\pi) + \mathbb{E} \left[\sum_t \gamma^t A_{\pi}(s_t, a_t) \right] \quad (\text{just rewriting})$$

$$= \eta(\pi) + \sum_s P_{\tilde{\pi}}(s) \sum_a \tilde{\pi}(a|s) A_{\pi}(s, a)$$

this is hard to do,
 so replace by δ_{π}

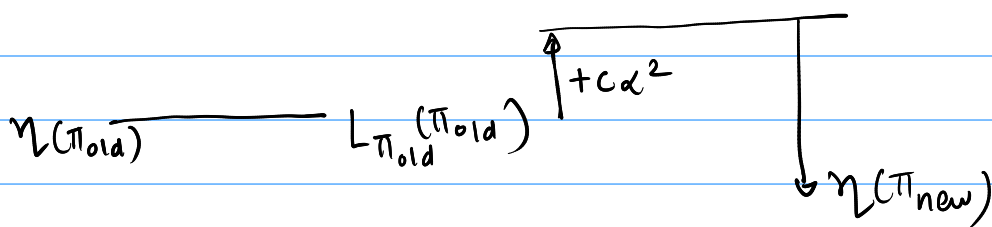
$$L_{\pi}(\tilde{\pi}) = \eta(\pi) + \sum_s \delta_{\pi}(s) \sum_a \tilde{\pi}(a|s) A_{\pi}(s, a)$$

Note: $L_{\pi}(\pi) = \eta(\pi) + 0$

Result in paper:

$$\eta(\pi_{\text{new}}) \leq L_{\pi_{\text{old}}}(\pi_{\text{new}}) + C\alpha^2$$

↗ constant
 ↘ α = mixture ratio



$$\eta(\tilde{\pi}) \leq L_{\pi}(\tilde{\pi}) + C D_{KL}^{\max}(\pi, \tilde{\pi})$$

Define $M_i(\pi) = L_{\pi_i}(\pi) + C D_{KL}^{\max}(\pi_i, \pi)$,

then $\eta(\pi_{i+1}) - \eta(\pi_i) \leq \underbrace{M_i(\pi_{i+1})}_{\text{minimise}} - \underbrace{M_i(\pi_i)}_{\text{equal to } \eta(\pi_i)}$

$\min_{\theta} [M_i(\pi)]$ is an unconstrained problem, but D_{KL}^{\max} is difficult
 \hookrightarrow make it a constrained version

$$\min_{\theta} L_{\theta_{old}}(\theta) \quad \text{subject to} \quad D_{KL}^{\max}(\theta_{old}, \theta) \leq \delta$$

still hard, so use sample average

After much ado:

$$\nabla L^T(\theta_{new} - \theta) \quad \min_{\theta} \mathbb{E} \left[\frac{\pi_{\theta}(a|s)}{q(a|s)} Q_{\theta_{old}}(s, a) \right]$$

replacement for advantage
sampling

$$\frac{1}{2}(\theta_{new} - \theta)^T F(\theta_{new} - \theta) \leq \delta \quad \text{subject to} \quad \mathbb{E}_{\theta_{old}} [D_{KL}(\pi_{\theta_{old}} \parallel \pi_{\theta})] \leq \delta \quad \frac{\partial \theta}{\partial \theta_i}$$

quad approx

How to solve?

\hookrightarrow first order approximation, no easy step size
 \hookrightarrow second order

\hookrightarrow need covariance matrix of gradients
 $n \times n$ sized
 \hookrightarrow instead, use conjugate gradients

Conjugate
gradients

$$\Delta \theta = \sqrt{\frac{2\delta}{\nabla L^T F(\theta)^{-1} \nabla L}} \begin{bmatrix} -1 \\ F \nabla L \end{bmatrix}$$

guess product directly

$$x = f(n)$$

$$x_1 = f(x_0) \\ x_2 = f(x_1)$$

$$x_{n+1} \approx f(x_n)$$

PPO

How can we avoid conjugate gradients and all the other drama?

Recall $L - D_{KL}^{\max}$ is the original objective,
converted to $\min L$ subject to $D_{KL} \leq \delta$

$$L(\theta) = \mathbb{E} \left[\underbrace{\frac{\pi_{\theta}(a|s)}{q(a|s)}}_{\text{use } \pi_{\theta_{old}} \text{ as sampler}} \underbrace{Q_{\theta}(s, a)}_{\text{revert to } A_t} \right]$$

conservative
policy
iteration

$$L^{CPI}(\theta) = \mathbb{E}_t \left[\underbrace{\frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}}_{r_t} \hat{A}_t \right]$$

$$L^{CLIP}(\theta) = \mathbb{E}_t \left[\min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1-\epsilon, 1+\epsilon) \hat{A}_t) \right]$$

We can even reintroduce D_{KL} in the objective
(see eqn 8 in paper), so long as its coefficient is
adaptively updated