

General Review of Terms

Policy: Dictates what action to take given a state.

Deterministic, $\Pi: s \rightarrow A$, maps states to actions

Stochastic, $\Pi(a|s)$, conditional probability distribution

MDP: 5 tuple $\{S(\text{set of states}),$

$A(\text{set of actions}),$

$P(s_{t+1}|s_t=s, a_t=a)$ aka $P(s'|s, a)$ (probability of state transition with action a),

$R(s'|s, a)$ (reward obtained from going to s' with a from s),

γ (discount factor, $\gamma=0$ is immediate reward, $\gamma=1$ is total reward)

Value Function: Value of state - expected reward of starting at state and continuing with policy

Written $V^\Pi(s), Q^\Pi(s, a)$, conditional to policy Π . If conditional to optimal policy, Π^* , then $V^*(s), Q^*(s, a)$.

$V^\Pi(s)$ expected value of following Π forever when agent starts at s .

$Q^\Pi(s, a)$ expected value of following Π forever when first taking action a in state s .

$$V^\Pi(s) = E_\Pi \{ R_t | s_t=s \} = E_\Pi \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t=s \right\}, \text{expected return starting from state given } \Pi$$

So, if $\gamma \in [0, 1] \ni 0, \sum_{k=0}^{\infty} \gamma^k r_{t+k} = r_{t+1}$, immediate reward. If $\gamma=1$, then r at all time steps equally weighted

Note, $R_t = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} \dots$

$$r_{t+1} + \gamma(r_{t+2} + \gamma r_{t+3} \dots) \quad E_\Pi \{ \cdot \} \text{ is expectation operator for policy } \Pi.$$

$$r_{t+1} + \gamma R_{t+1} \Rightarrow$$

$$V^\Pi(s) = E_\Pi \{ R_t | s_t=s \}$$

$$E_\Pi \{ r_{t+1} + \gamma V^\Pi(s_{t+1}) | s_t=s \}$$

Policy Gradient

$$\text{Objective function } J(\theta) = E_\theta \left[\sum_{t=0}^{T-1} R(s_t, a_t; \theta) \right] = \sum_{t=0}^{T-1} P(r; \theta) R(t) \quad \leftarrow \text{stochastic policy}$$

Find policy θ that creates trajectory T $(s_0, a_0, s_1, a_1, s_2, a_2, \dots, s_{T-1}, a_{T-1})$

that maximizes expected return, max $J(\theta) = \max_{\theta} \sum_{t=0}^{T-1} P(r; \theta) R(t)$

Expectation in continuous state space $E_{p(s|s_0)}[J(\theta)] = \int p(s) J(\theta) ds$

Derivations

$$\text{Note: } \nabla_\theta J(\theta) = J(\theta) \frac{\partial J(\theta)}{\partial \theta} = J(\theta) \nabla_\theta \ln J(\theta)$$

$$\text{Let } J(\theta) = E_{\pi_\theta(T)}[r(T)] = \int \pi_\theta(r) r(T) dT, \quad \pi_\theta(r) = P(r) \text{ given } \Pi_\theta \quad \leftarrow \text{Expectation of } T \text{ is Prob}(T) \cdot \text{reward}(T),$$

$$\text{Now } \nabla_\theta J(\theta) = \int \nabla_\theta \pi_\theta(r) r(T) dT$$

$$= \int \pi_\theta(r) \underbrace{\nabla_\theta \ln \pi_\theta(r) r(T)}_{\text{let this be } \nabla_\theta r(T)} dT$$

$$= E_{\pi_\theta(T)}[r(T) \nabla_\theta \ln \pi_\theta(r)]$$

Thus, policy gradient is an expectation.

$$\nabla_\theta \ln \pi_\theta(T) \rightarrow \text{First, find } \pi_\theta(T) = \pi(s_0, a_0, s_1, a_1, \dots, s_{T-1}, a_{T-1}) = P(s_0) \prod_{t=1}^{T-1} \pi_\theta(a_t|s_t) P(s_{t+1}|s_t, a_t)$$

$$\text{Then, } \ln \pi_\theta(T) = \ln(P(s_0) \prod_{t=1}^{T-1} \pi_\theta(a_t|s_t) P(s_{t+1}|s_t, a_t))$$

$$= \ln P(s_0) + \sum_t \ln \pi_\theta(a_t|s_t) + \sum_{t=1}^{T-1} \ln P(s_{t+1}|s_t, a_t)$$

Now, apply ∇_θ

$$\nabla_\theta \ln P(s_0) + \sum_t \nabla_\theta \ln \pi_\theta(a_t|s_t) + \sum_{t=1}^{T-1} \nabla_\theta \ln P(s_{t+1}|s_t, a_t)$$

$$\nabla_\theta \ln \pi_\theta(T) = \sum_t \nabla_\theta \ln \pi_\theta(a_t|s_t)$$

$$\nabla_\theta J(\theta) = E_{\pi_\theta(T)} \left[r(T) \sum_t \nabla_\theta \ln \pi_\theta(a_t|s_t) \right] \quad \leftarrow \text{from above}$$

$$\nabla_\theta J(\theta) = \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_\theta \ln \pi_\theta(a_{i,t}|s_{i,t}) \right) \left(\sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right)$$

$$\left. \begin{array}{c} \text{within a } T, \text{ this is } \nabla_\theta \text{ of expectation} \\ \uparrow \qquad \qquad \qquad \uparrow \\ \text{over all } T \qquad \text{likelihood of } T \qquad \text{reward of } T \end{array} \right)$$

$$\text{Remember: } E \left[\sum_{t=0}^{T-1} R(s_t, a_t; \theta) \right] = \sum_{t=0}^{T-1} P(r; \theta) R(t)$$

$$\left. \begin{array}{c} \uparrow \\ \text{for arbitrary states} \\ \text{actions given policy} \\ \text{For } t \in \{0, \dots, T-1\}, \text{ over } T \end{array} \right)$$

Can maximize rewards of θ over expectation pertaining to π

$$\Theta = \Theta + \alpha \nabla_{\theta} J$$

Note: $\sum_{t=1}^T \nabla_{\theta} \ln \Pi_{\theta}(a_{i,t} | s_{i,t})$ is max log likelihood from $\nabla_{\theta} \ln \pi(a) = \pi(a) \nabla_{\theta} \ln \pi(a)$
 Gradient term of $\prod_{t=1}^T \Pi_{\theta}(a_t | s_t)$ from missing grad (summand uses instead of \prod).

Now for Cart-Pole-v1...

$$J_{\theta}(\pi) = \sum_{i=1}^{N_s} \left(\sum_{t=1}^T \nabla_{\theta} \ln \Pi_{\theta}(a_{i,t} | s_{i,t}) \right) \left(\sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right)$$

could use logistic since $|A|=2$. But softmax generalizes.

So, need: parametrized policy which tells us probability of action given state

$$\lambda \in \{0, 1\}, \text{ let } \Pi(a_i | s), \Pi(a_i | s) = \text{softmax with } \Pi(a_i | s) = \frac{e^{w_i s}}{\sum_{j=0}^{|A|-1} e^{w_j s}}$$

Training to derive: $\frac{\partial J(\pi)}{\partial w_{ij}}$

$$\begin{matrix} & 2 \times 4 & 4 \times 1 \\ \{ w_i \rightarrow \boxed{\quad \quad \quad \quad} \cdot \underbrace{\boxed{s}}_{\substack{2 \times 1}} = \boxed{s_i} \rightarrow \text{softmax}(\boxed{s_i}) = \boxed{\frac{P(a_i | s)}{P(a_i | s)}} \end{matrix}$$

Update rule for single weight $w_{ij} = i$ row, column

of system.

$\frac{\partial}{\partial w_{ij}} \ln \Pi(a_i | s) \leftarrow \Pi$ dependent on $i \neq j$, due to bottom of softmax term. Could use Kronecker delta, implement with one-hot encoding,

$$\begin{aligned} \frac{\partial}{\partial w_{ij}} (\ln e^{w_i s} - \ln \sum_{j=0}^{|A|-1} e^{w_j s}) \text{ for } i=j &\quad \text{Note: } i=\text{action}, j=\text{state dimension} \quad \text{For } i \neq j \\ s_j - \left(\sum_{j=0}^{|A|-1} e^{w_j s} \right)^{-1} \frac{\partial}{\partial w_{ij}} \left(\sum_{j=0}^{|A|-1} e^{w_j s} \right) &\quad \frac{\partial}{\partial w_{ij}} (\ln e^{w_i s} - \ln \sum_{j=0}^{|A|-1} e^{w_j s}) \\ s_j - \left(\sum_{j=0}^{|A|-1} e^{w_j s} \right)^{-1} s_j e^{w_i s} &\quad - \left(\sum_{j=0}^{|A|-1} e^{w_j s} \right)^{-1} s_j e^{w_i s} \\ s_j - s_j \Pi(a_i | s) &\quad - s_j \Pi(a_i | s) \\ s_j (1 - \Pi(a_i | s)) &\quad \text{Thus, } \frac{\partial}{\partial w_{ij}} \ln \Pi(a_i | s) = \begin{cases} s_j (1 - \Pi(a_i | s)) & \text{if } i=j \\ -s_j \Pi(a_i | s) & \text{if } i \neq j \end{cases} \end{aligned}$$

$$\Delta w_{ij} \in \Pi_w, \quad \frac{\partial}{\partial w_{ij}} \ln \Pi(a_i | s) = s_j (\delta - \Pi(a_i | s)), \text{ where } \begin{cases} \delta = 1 & \text{if } i=j \\ \delta = 0 & \text{if } i \neq j \end{cases}$$

Update Rule: $w_{ij} = w_{ij} + \alpha \frac{\partial J}{\partial w_{ij}}$

for single π

$$= w_{ij} + \alpha \sum_t [s_j (\delta - \Pi(a_i | s_t)) r(s_t, a_t)]$$

Vectorized: $w = w + \alpha \sum_t [\delta - \Pi(a | s_t) r(s_t, a_t)] s^T$, where δ is one-hot of action taken at step t at π

Given $w \in \mathbb{R}^{|A|} \times \mathbb{R}^{1 \times |S|}$, $s \in \mathbb{R}^{1 \times |S|}$, $\Pi(a | s) = \mathbb{R}^{|A|} \times 1$, where s_t is features of s .
 softmax parameter matrix state input output of softmax

Problem: $w = w - \alpha \sum_t [\delta - \Pi(a | s_t) r(s_t, a_t)] s^T$ seems to converge (subtracting gradient)!?

$(\delta - \Pi(a | s_t))$ should say how far our policy was from action

\pm existing weights to get to that action

$r(s_t, a_t)$ tells us how much we want to get to that action

if high, amplifies weight change, if low, reduces weight change

Adding this to weights increases prob. of high reward choices... so why does subtracting work?!