

Proximal Policy Optimization

Reinforcement Learning

November 17, 2020

MSU

Reminder: Lower Bound

$$J_{p,q} - J_{p^{\text{old}},q} \geq L_{\text{old}}(p,q):$$

Reminder: Lower Bound

importance sampling

$$J_{p,q} - J_{p^{old},q} \leq L_{old,p,q} \left(\frac{1}{n} \sum_{i=1}^n \frac{p(x_i) | q(x_i)}{p^{old}(x_i) | q(x_i)} - 1 \right)$$

data generated by p^{old}

do not require fresh critic

Reminder: Lower Bound

importance sampling

$$J_{p,q} - J_{p^{old},q} \leq L_{old,p,q} \left(\frac{1}{n} \sum_{i=1}^n \frac{p(x_i) | q(x_i)}{p^{old}(x_i) | q(x_i)} - 1 \right)$$

data generated by p^{old}

do not require fresh critic

Approximation Error Bound:

$$C K L_{p^{old},k,q}$$

Reminder: Lower Bound

importance sampling

$$L_{old}(p, q) = \frac{1}{n} \sum_{i=1}^n \frac{p(x_i)}{q(x_i)} \log \frac{p(x_i)}{q(x_i)}$$

data generated by q

do not require fresh critic

Approximation Error Bound:
 $C \sqrt{\frac{KL(p, q)}{n}}$

$$L_{old}(p, q) - C \sqrt{\frac{KL(p, q)}{n}} \leq L(p, q) \leq L_{old}(p, q) + C \sqrt{\frac{KL(p, q)}{n}}$$

Reminder: Lower Bound

importance sampling

$$J_{p,q} - J_{p^{old},q} \approx L_{p,q} - L_{p^{old},q} = \frac{1}{N} \sum_{i=1}^N \frac{p(x_i) - p^{old}(x_i)}{p^{old}(x_i)} \log \frac{p(x_i)}{p^{old}(x_i)}$$

data generated by p^{old}

do not require fresh critic

Approximation Error Bound:
 $C \sqrt{\frac{KL(p^{old} || q)}{N}}$

$$L_{p,q} - L_{p^{old},q} \leq C \sqrt{\frac{KL(p^{old} || q)}{N}}$$

can't compute constant C ;
 theoretically it is huge;
 critic is imperfect;

Reminder: Trust Region Policy Optimization (TRPO)

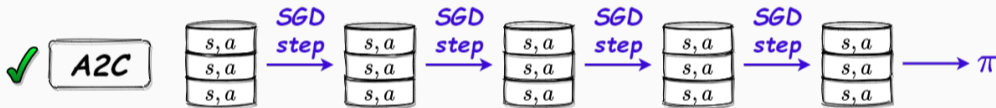
$$\begin{aligned} & \max_{\theta} L(\theta; p, q, \tilde{N}) \\ & \text{s.t. } \text{KL}(p \parallel q) \leq \alpha \end{aligned}$$

X robust: prevents large changes;

Reminder: Trust Region Policy Optimization (TRPO)

$$\begin{aligned} & \text{maximize } \mathbb{E}_{p \sim q} L_{\text{old}} \\ & \text{subject to } \text{KL}(p \text{ old} \| k) \leq \alpha \end{aligned}$$

X robust: prevents large changes;



Reminder: Trust Region Policy Optimization (TRPO)

$$\begin{aligned} & \max_{\theta} L(\theta; p, q, \tilde{N}) \\ & \text{s.t. } \text{KL}(p_{\theta} \parallel p_{\text{old}}) \leq \alpha \end{aligned}$$

X robust: prevents large changes;

Reminder: Trust Region Policy Optimization (TRPO)

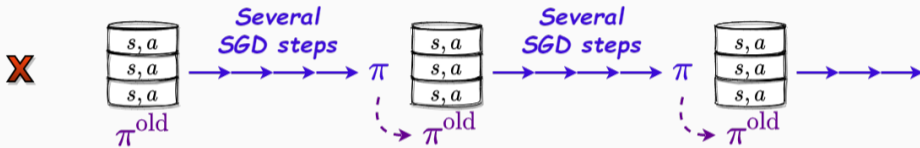
$$\begin{aligned} & \max_{\theta} L(\theta; p, q, \tilde{N}) \\ & \text{s.t. } \text{KL}(p_{\theta} \parallel p_{\text{old}}) \leq \alpha \end{aligned}$$

X robust: prevents large changes;

critic and actor can't share backbone;
computationally costly;
complicated :(

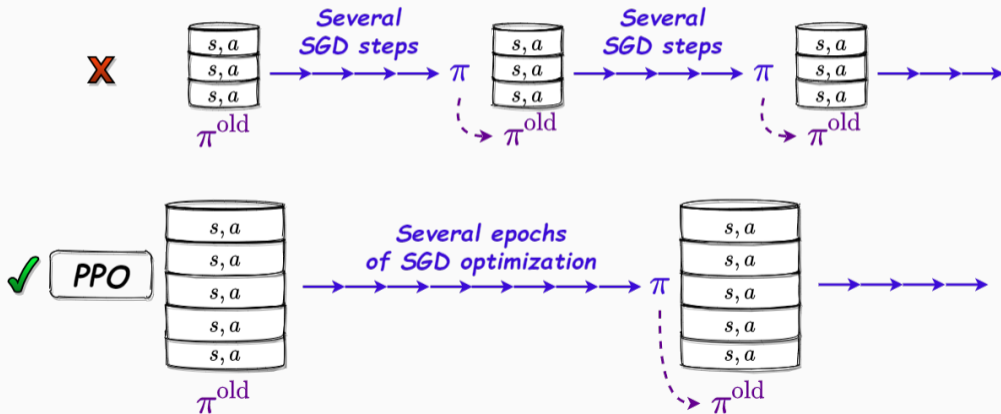
Proximal Policy Optimization (PPO): Pipeline

$$E_s \left[d_{\pi^{\text{old}}} \right] E_a \left[\frac{p(a|s)}{p^{\text{old}}(a|s)} A^{\text{old}} \right] \text{C K L p}^{\text{old}} k \quad q \tilde{N} \max$$



Proximal Policy Optimization (PPO): Pipeline

$$E_s \left[d_{\pi^{\text{old}}} \right] E_a \left[\frac{p(a|s)}{p_{\text{old}}(a|s)} A \right] \text{C KLp}^{\text{old}} k q \tilde{N} \max$$



Clipping Objective

$$L_{\text{old}}(p, q) = \text{C}_{\text{KL}}(p^{\text{old}} \| k) \quad q \sim \tilde{N} \quad \max$$

Default surrogate function:

$$p, q: \frac{p(a) | s(q)}{\text{old } p(a) | s(q)}$$

$$L_{\text{old}}(p, q) = \mathbb{E}_{s; a} [p(q) A^{\text{old}}(p; a) q]$$

Clipping Objective

$$L_{\text{old}}(p, q) = \mathbb{E}_{\mathcal{K}} \text{KL}(p^{\text{old}} \| q) \quad \tilde{N} \text{ max}$$

Default surrogate function:

$$p \llbracket q : \frac{p \wedge q}{p \vee q}$$

$$L_{\text{old}}(p, q) = \mathbb{E}_{s; a} p \wedge q \wedge A^{\text{old}}(p; a) \vee q$$

Clipping Objective

$$L_{\text{old}}(p, q) = C \text{KL}(p^{\text{old}} \| q) \tilde{N} \max$$

Default surrogate function:

$$p, q: \frac{pa | sq}{\text{old} pa | sq}$$

$$L_{\text{old}}(p, q) = E_{s; a} p, q A^{\text{old}} ps; aq$$

Clipped surrogate function:

$$\text{clip} p, q: \text{clipp } p, q; 1, 1, q$$

$$L_{\text{old}}^{\text{clip}}(p, q) = E_{s; a} \text{clip } p, q A^{\text{old}} ps; aq$$

Clipping Objective

$$L_{\text{old}}(p, q) = C \text{KL}(p^{\text{old}} \| q) \tilde{N} \max$$

Default surrogate function:

$$p, q: \frac{pa | sq}{\text{old} pa | sq}$$

$$L_{\text{old}}(p, q) = E_{s; a} p, q A^{\text{old}} ps; aq$$

Clipped surrogate function:

$$\text{clip} p, q: \text{clipp } p, q; 1 \quad ; 1 \quad q$$

$$L_{\text{old}}^{\text{clip}}(p, q) = E_{s; a} \text{clip } p, q A^{\text{old}} ps; aq$$

Recalling lower bound intuition

$$\mathbb{E}_s \mathbb{E}_a \min_p \left[\underbrace{d_{old}(p|s)}_{\text{original term}} + \underbrace{\lambda \sum_{a \in \mathcal{A}} \underbrace{w_a}_{\text{importance sampling weight}} \underbrace{d_{old}(a|p)}_{\text{regularization}}}_{\text{term with clipped importance sampling weight}} \right]$$

The diagram illustrates the decomposition of the lower bound expression. It shows the expectation over states s and actions a of the minimum over policies p of a sum of two terms. The first term is the original KL divergence $d_{old}(p|s)$. The second term is a regularization term $\lambda \sum_{a \in \mathcal{A}} w_a d_{old}(a|p)$, where w_a is the importance sampling weight. The regularization term is shown to be equivalent to a clipped version $\text{clip}(w_a, C, K)$.

Recalling lower bound intuition

$E_s d_{old} p_{sq} E_a d_{old} p_{sq} \min_p p q A^{old} p s; a q; \text{clip} p q A^{old} p s; a q q C K L p^{old} k q \tilde{N} \max$

original term importance sampling weight regularization

Advantage Sign	Direction	Bad ratio case	Gradient
$A^{old} p s; a q \neq 0$			

Recalling lower bound intuition

E_s d_{old} p_{sq} E_a $d_{pa|sq}$ \min_p p q A^{old} $ps; aq;$ clip p q A^{old} $ps; aq$ q C K L p k q \tilde{N} \max

original term importance sampling weight regularization

Advantage Sign	Direction	Bad ratio case	Gradient
$A^{old} ps; aq \neq 0$	$pa sq \partial$		

Recalling lower bound intuition

$E_s d_{old} p_{sq} E_a \quad \text{original term} \quad \text{importance sampling weight} \quad \text{term with clipped} \quad \text{regularization}$
 $\min_p p q A^{old} p s; a q; \quad \text{clip} \quad p q A^{old} p s; a q \quad q \quad C K L p^{old} k \quad q \tilde{N} \max$

Advantage Sign	Direction	Bad ratio case	Gradient
$A^{old} p s; a q \neq 0$	$p a s q \partial$	$p q i \quad 1:2$	

Recalling lower bound intuition

$E_s d_{old} p_{sq} E_a \quad \text{original term} \quad \text{importance sampling weight} \quad \text{term with clipped} \quad \text{regularization}$
 $\min_p p q A^{old} p s; a q; \quad \text{clip} \quad p q A^{old} p s; a q \quad q \quad C K L p^{old} k \quad q \tilde{N} \max$

Advantage Sign	Direction	Bad ratio case	Gradient
$A^{old} p s; a q \neq 0$	$p a s q \partial$	$p q i \quad 1:2$	0

Recalling lower bound intuition

$$E_{s \sim p} \left[\sum_{i=1}^N \left(\frac{1}{\sqrt{p_i}} \left(\frac{p_i}{q_i} \right)^{\frac{1}{2}} \right) \right] \geq \frac{1}{\sqrt{2}} \sqrt{\sum_{i=1}^N \frac{p_i}{q_i}}$$

original term: $\frac{1}{\sqrt{p_i}} \left(\frac{p_i}{q_i} \right)^{\frac{1}{2}}$
 importance sampling weight: $\frac{1}{\sqrt{p_i}}$
 regularization: $\left(\frac{p_i}{q_i} \right)^{\frac{1}{2}}$
 clip: $\frac{1}{\sqrt{2}}$
 $\min(p_i, q_i)$

Advantage Sign	Direction	Bad ratio case	Gradient
$A^{\text{old}} \geq 0$	$p \leq q$	$p/q \in [1, 2]$ $p/q \in [0, 8]$	0

Recalling lower bound intuition

$$E_{s \sim p} \left[\sum_{i=1}^N \left(\frac{1}{\sqrt{p_i}} \log \frac{p_i}{q_i} \right) \right] \geq \sum_{i=1}^N \left(\frac{1}{\sqrt{p_i}} \log \frac{p_i}{q_i} \right) - \frac{1}{2} \sum_{i=1}^N \frac{1}{p_i} \log \frac{p_i}{q_i} - \frac{1}{2} \sum_{i=1}^N \frac{1}{q_i} \log \frac{p_i}{q_i}$$

original term
importance sampling weight
regularization

Advantage Sign	Direction	Bad ratio case	Gradient
$A^{old} p_i; a_i \neq 0$	$p_i s_i \leq 0$	$p_i q_i \leq 1:2$ $p_i q_i \geq 0:8$	0 same

Recalling lower bound intuition

$$E_{s \sim p} \left[\sum_{i=1}^N \frac{1}{\sqrt{p_i}} \left(\frac{A_i}{\sqrt{p_i}} - \frac{A_i^{\text{old}}}{\sqrt{p_i}} \right)^2 \right] \leq \frac{1}{\epsilon} \left(\sum_{i=1}^N \frac{1}{\sqrt{p_i}} \left(\frac{A_i}{\sqrt{p_i}} - \frac{A_i^{\text{old}}}{\sqrt{p_i}} \right)^2 + \sum_{i=1}^N \frac{1}{\sqrt{p_i}} \left(\frac{A_i}{\sqrt{p_i}} - \frac{A_i^{\text{old}}}{\sqrt{p_i}} \right)^2 \right)$$

original term importance sampling weight regularization
 term with clipped

Advantage Sign	Direction	Bad ratio case	Gradient
$A^{\text{old}} \neq 0$	$p \mid q \rightarrow$	$p \mid q \quad 1:2$ $p \mid q \quad 0:8$	0 same
$A^{\text{old}} = 0$	$p \mid q \leftarrow$	$p \mid q \quad 1:2$ $p \mid q \quad 0:8$	

Recalling lower bound intuition

$$E_{s \sim p} \left[\sum_{i=1}^N \frac{1}{\sqrt{p_i}} \left(\frac{A_{i,old}}{p_i} - \frac{A_{i,clip}}{p_i} \right) \right] \leq \frac{C}{\sqrt{N}}$$

original term
importance sampling weight
regularization

Advantage Sign	Direction	Bad ratio case	Gradient
$A_{old} \neq 0$	$p_a > p_b$	$p_a : p_b = 1:2$ $p_a : p_b = 0:8$	0 same
$A_{old} = 0$	$p_a > p_b$	$p_a : p_b = 1:2$ $p_a : p_b = 0:8$	same 0

Recalling lower bound intuition

$$E_{s \sim p} [d_{old}(p, s)] - E_{a \sim p} [d_{old}(p, a)] \min_{p, q} p \cdot q \cdot A^{old}(p; a, q) \cdot \text{clip} \left(\frac{p \cdot q \cdot A^{old}(p; a, q)}{C \cdot K \cdot L \cdot p^{old} \cdot k}, q \right) \tilde{N} \max$$

original term importance sampling weight regularization
term with clipped

Advantage Sign	Direction	Bad ratio case	Gradient
$A^{old}(p; a, q) \neq 0$	$p \cdot a \mid s \cdot q \cdot \tilde{N}$	$p \cdot q \cdot i \quad 1:2$ $p \cdot q \quad 0:8$	0 same
$A^{old}(p; a, q) = 0$	$p \cdot a \mid s \cdot q \cdot \tilde{N}$	$p \cdot q \cdot i \quad 1:2$ $p \cdot q \quad 0:8$	same 0

Clipped Critic Loss

$$\text{Loss}_\phi = \mathbb{E}_q [p(y) \min(V, p(q)^2)]$$

Clipped Critic Loss

$$\text{Loss}_q: p(y) \quad V(p, q) \\ p(y) \quad V^{\text{old}} \quad V^{\text{old}} \quad V(p, q)$$

Clipped Critic Loss

$$\text{Loss}_{\text{critic}} = \mathbb{E}_{p(y)} [V(p, q)] - \mathbb{E}_{p(y)} [V^{\text{old}}(p, q)]$$

$$\text{Loss}_{\text{clip}} = \mathbb{E}_{p(y)} [V^{\text{old}}(p, q) - \text{clip}(V^{\text{old}}(p, q), V(p, q), \epsilon)]$$

Clipped Critic Loss

$$\text{Loss}_p q: p y - V p q^2$$

$$p y - V^{\text{old}} - V^{\text{old}} - V p q^2$$

$$\text{Loss}^{\text{clip}} p q: p y - V^{\text{old}} - \text{clip}(V^{\text{old}} - V p q, \epsilon) - V p q^2$$

$$\max_p \text{Loss}_p q - \text{Loss}^{\text{clip}} p q$$

Bias-Variance trade-o

Given rollout $s; r; s^1; r^1; s^2; r^2 \dots s^{pMq}$ from policy and approximation of $V_{\psi q}$

Bias-Variance trade-o

Given rollout $s; r; s^1; r^1; s^2; r^2 \dots s^{M_q}$ from policy π and approximation of V_{π}
perform **credit assingment** for state-action pairs; a (was this decision good or bad?)

Bias-Variance trade-o

Given rollout $s; r; s^1; r^1; s^2; r^2 \dots s^{PMq}$ from policy and approximation of V_{π}
perform **credit assingment** for state-action pairs; a (was this decision good or bad?)

For Actor:

$r : p \quad q \quad \log \quad p \mid s, a$
look on
advantage
estimator

For Critic:

V_{π}
look on
target
for regression

Bias-Variance trade-o

Given rollout $s; r; s^1; r^1; s^2; r^2 \dots s^{Mq}$ from policy π and approximation of V_{π}
 perform **credit assingment** for state-action pairs; a (was this decision good or bad?)

For Actor:

For Critic:

$r : \pi(a|s) \log \frac{\pi(a|s)}{p(a|s)}$
 advantage estimator

$V_{\pi} = \mathbb{E}_{\pi} [r^2 | s]$
 target for regression

	$\pi(a s)$	Bias	Variance
Monte Carlo	$\mathbb{E}_{\pi} [r^2 s]$	0	high
1-step	$r - V_{\pi}(s)$	high	low

Bias-Variance trade-o

Given rollout $s; r; s^1; r^1; s^2; r^2 \dots s^{Mq}$ from policy π and approximation of $V_{\pi; aq}$ perform **credit assingment** for state-action pairs; a (was this decision good or bad?)

For Actor:

For Critic:

$r : \pi(a|s) \log \frac{\pi(a|s)}{\pi(a|s)}$
 advantage estimator

$V_{\pi; aq}$
 target for regression

	$\pi; aq$	Bias	Variance
Monte Carlo	$\pi; aq: r + \gamma V_{\pi; aq}$	0	high
N-step	$\pi; aq: r + \gamma V_{\pi; aq}$	intermediate	intermediate
1-step	$\pi; aq: r + \gamma V_{\pi; aq}$	high	low

Bias-Variance trade-o

Given rollout $s; r; s^1; r^1; s^2; r^2 \dots s^{Mq}$ from policy π and approximation of $V_{\pi}(s)$ perform **credit assignment** for state-action pairs; a (was this decision good or bad?)

For Actor:

For Critic:

$r : \pi(a|s) \log \frac{\pi(a|s)}{\pi(a|s)}$
 advantage estimator

$V_{\pi}(s) : \pi(a|s) V_{\pi}(s)$
 target for regression

	$\pi(a s)$	Bias	Variance
Monte Carlo	$\pi(a s) : r + \gamma V_{\pi}(s) - V_{\pi}(s)$	0	high
N-step	$\pi(a s) : r + \gamma V_{\pi}(s) - V_{\pi}(s)$	intermediate	intermediate
1-step	$\pi(a s) : r + \gamma V_{\pi}(s) - V_{\pi}(s)$	high	low

Problem: hard to choose N .

Backward view: idea

N-step update:

$$V_{psq} \leftarrow V_{psq} + \rho_N (q_{ps} - a_{psq})$$

Backward view: idea

N-step update:

$$V(p,s,q) \approx V(p,s,q) + \gamma [r + V(p,s',q) - V(p,s,q)]$$

How to turn 1-step update into 2-step?

$$V(p,s,q) \approx V(p,s,q) + \gamma [r + V(p,s',q) - V(p,s,q)] + \gamma^2 [r + V(p,s',q) - V(p,s,q)]$$

Backward view: idea

N-step update:

$$V_{psq} \ominus V_{psq} = \sum_{a \in \mathcal{A}} \rho^N Q_{ps,a} - Q_{ps}$$

How to turn 1-step update into 2-step?

$$V_{psq} \ominus V_{psq} = r + \gamma V_{ps^1q} - V_{psq} + \gamma V_{ps^2q} - V_{ps^1q} + V_{ps^1q} - \gamma V_{ps^2q} + V_{psq} - \gamma V_{ps^1q}$$

Backward view: idea

N-step update:

$$V_{psq} \leftarrow V_{psq} + \rho \sum_{a \in \mathcal{A}} \sum_{s'} \gamma^t \left(V_{ps'a} - V_{psq} \right) \pi(a|s)$$

How to turn 1-step update into 2-step?

$$V_{psq} \leftarrow V_{psq} + \rho \left(r + \gamma V_{ps'a} - V_{psq} \right) \pi(a|s)$$

Backward view: idea

N-step update:

$$V_{psq} \leftarrow V_{psq} + \rho \sum_{q'} \Delta_{psq} \Delta_{psq}$$

How to turn 1-step update into 2-step?

$$\begin{array}{c}
 V_{psq} \leftarrow V_{psq} + \rho \sum_{q'} \Delta_{psq} \Delta_{psq} \\
 \\
 V_{psq} \leftarrow V_{psq} + \rho \sum_{q'} \Delta_{psq} \Delta_{psq} + \rho \sum_{q'} \Delta_{psq} \Delta_{ps^1q} \\
 \\
 V_{psq} \leftarrow V_{psq} + \rho \sum_{q'} \Delta_{psq} \Delta_{psq} + \rho \sum_{q'} \Delta_{psq} \Delta_{ps^1q} + \rho \sum_{q'} \Delta_{psq} \Delta_{ps^2q} + \dots + \rho \sum_{q'} \Delta_{psq} \Delta_{ps^1q}
 \end{array}$$

Backward view: idea

N-step update:

$$V(p, s, q) \leftarrow V(p, s, q) + \alpha [r + \sum_{s'} P(s'|s, a) V(p, s', q) - V(p, s, q)]$$

How to turn 1-step update into 2-step?

$$V(p, s, q) \leftarrow \sum_{t=0}^{N-1} \gamma^t [r + \sum_{s'} P(s'|s, a) V(p, s', q)] + \gamma^N V(p, s, q)$$

$$V(p, s, q) \leftarrow V(p, s, q) + \alpha [r + \sum_{s'} P(s'|s, a) V(p, s', q) - V(p, s, q)] + \gamma [V(p, s, q) - V(p, s, q)]$$

Eligibility Traces

Eligibility Traces

Use 1-step TD-error to update $V_{\psi(s)}$ for all states

Eligibility Traces

Use 1-step TD-error to update $V(s)$ for all states

Define **eligibility trace** $e_t(s)$ as a coefficient of update:

$$\delta_t = V(s_{t+1}) - V(s_t) \quad e_t(s) = \rho e_{t-1}(s) + \mathbb{1}_{s_t=s}$$

Eligibility Traces

Use 1-step TD-error to update $V_{\psi(s)}$ for all states

Define **eligibility trace** $e_{\psi(s)}$ as a coefficient of update:

$$\delta_{\psi(s)}: V_{\psi(s)} \leftarrow V_{\psi(s)} + \delta_{\psi(s)} e_{\psi(s)}$$

Online Monte-Carlo updates:

$e_{\psi(s)} = 0$ at the start of each episode

Eligibility Traces

Use 1-step TD-error to update $V_{\psi(s)}$ for all states

Define **eligibility trace** $e_{\psi(s)}$ as a coefficient of update:

$$\Delta_{\psi(s)} = V_{\psi(s)} - V_{\psi(s)} + e_{\psi(s)} \delta_{\psi(s)}$$

Online Monte-Carlo updates:

- $e_{\psi(s)} = 0$ at the start of each episode
- $e_{\psi(s)} \leftarrow \gamma e_{\psi(s)} + 1$ after visitings

Eligibility Traces

Use 1-step TD-error to update V_{psq} for all states

Define **eligibility trace** e_{psq} as a coefficient of update:

$$\delta_{psq} = V_{psq} - V_{psq} \quad e_{psq} = \rho_{psq}$$

Online Monte-Carlo updates:

- $\delta_{psq} = 0$ at the start of each episode
- $e_{psq} \leftarrow \rho_{psq} + \gamma e_{psq}$ after visitings
- $\delta_{psq} \leftarrow \rho_{psq} - \gamma e_{psq}$ after each step

TD(1) and TD(0)

TD(1)

Input: policy

Initialize V arbitrarily

Initialize $\epsilon = 0$

observe s_0

for $k = 0; 1; 2; \dots$:

\hat{a}_k take action a_k , observe $r_k; s_{k+1}$

TD(1) and TD(0)

TD(1)

Input: policy

Initialize V ψ arbitrarily

Initialize $\epsilon = 0$

observe s_0

for $k = 0; 1; 2; \dots$:

\hat{a} take action a_k , observe $r_k; s_{k+1}$

\hat{p} $\psi_{k+1} = r_k + \gamma \psi_k - V_{k-1} + V_k$

TD(1) and TD(0)

TD(1)

Input: policy

Initialize V ψ arbitrarily

Initialize $\epsilon = 0$

observe s_0

for $k = 0; 1; 2; \dots$:

\hat{a} take action a_k , observe $r_k; s_{k+1}$

\hat{p} $\psi_{k+1} = r_k + \gamma \psi_k$

$\hat{\epsilon} = \epsilon + \beta$

TD(1) and TD(0)

TD(1)

Input: policy

Initialize V ψ arbitrarily

Initialize $\epsilon = 0$

observe s_0

for $k = 0; 1; 2; \dots$:

\hat{a} take action a_k , observe $r_k; s_{k+1}$

\hat{p} $\delta = r_k + V(\psi_{k+1}) - V(\psi_k)$

$\hat{\epsilon} = \epsilon + \delta$

$\hat{\psi}$: $V(\psi) \leftarrow V(\psi) + \epsilon \delta \psi$

TD(1) and TD(0)

TD(1)

Input: policy

Initialize V ψ arbitrarily

Initialize ϵ 0

observe s_0

for $k = 0; 1; 2; \dots$:

\hat{a} take action a_k , observe $r_k; s_{k+1}$

\hat{p} $\delta = r_k + V(\psi_{k+1}) - V(\psi_k)$

\hat{e} $\psi_{k+1} \leftarrow \psi_k + \delta \psi_k$

\hat{e} $V \leftarrow V + \epsilon \delta$

\hat{e} $\epsilon \leftarrow \epsilon \gamma$

TD(1) and TD(0)

TD(1)

Input: policy

Initialize V $p(s)$ arbitrarily

Initialize $\epsilon = 0$

observe s_0

for $k = 0; 1; 2; \dots$:

\hat{a} take action a_k , observe $r_k; s_{k+1}$

$\hat{p}(s)$: $r_k + V(p_{k+1}(s)) - V(p_k(s))$

$\hat{\epsilon} = \epsilon + \delta$

\hat{V} : $V(p(s)) \leftarrow V(p(s)) + \epsilon \hat{p}(s)$

$\hat{\epsilon}$: $\epsilon \leftarrow \epsilon - \epsilon$

TD(0)

Input: policy

Initialize V $p(s)$ arbitrarily

Initialize $\epsilon = 0$

observe s_0

for $k = 0; 1; 2; \dots$:

\hat{a} take action a_k , observe $r_k; s_{k+1}$

$\hat{p}(s)$: $r_k + V(p_{k+1}(s)) - V(p_k(s))$

$\hat{\epsilon} = \epsilon + \delta$

\hat{V} : $V(p(s)) \leftarrow V(p(s)) + \epsilon \hat{p}(s)$

$\hat{\epsilon}$: $\epsilon \leftarrow \epsilon - \epsilon$

TD()

TD()

TD()

TD()

TD()

Input: policy

Initialize V arbitrarily

Initialize $\epsilon = 0$

observe s_0

for $k = 0; 1; 2; \dots$:

\hat{a}_k take action a_k , observe $r_k; s_{k+1}$

$\hat{p}_{k+1} = r_k + \gamma V(s_{k+1}) - V(s_k)$

$\epsilon_{k+1} \in [0, 1]$

$\hat{V}_{k+1} = V(s_k) + \epsilon_{k+1} \hat{p}_{k+1}$

$\epsilon_{k+1} \in [0, \epsilon_k]$

Backward view vs Forward view

Forward View

Give credit to present from known future

is this decision good or bad based on the
outcome?

Backward view vs Forward view

Forward View

Give credit to present from known future

is this decision good or bad based on the outcome?

Backward View

Update past credits with present information

which decisions in the past to blame?

Forward view for TD(γ)

Equivalent forms of TD(γ) updates

$$s_t \rightarrow s_{t+1} \quad p_t \rightarrow p_{t+1} \quad a_t \rightarrow a_{t+1}$$
$$V_t \rightarrow V_{t+1} \quad \delta = r_t + \gamma V_{t+1} - V_t$$

Forward view for TD(γ)

Equivalent forms of TD(γ) updates

$$p_t = \sum_{s=0}^{\infty} \gamma^s p_{t+s} \quad ; \quad a_t = \sum_{s=0}^{\infty} \gamma^s a_{t+s}$$

Forward view for TD(γ)

Equivalent forms of TD(γ) updates

$$p_t^s + \gamma (r_t + p_t^s - \sum_{q=1}^N p_{t-1}^q w_{s,q}^{(t)} - \sum_{q=1}^N p_{t-1}^q w_{s,q}^{(t)})$$

Term	Left side	Left side coe .	Right side	Right side coe .
r_t^s				

Forward view for TD(γ)

Equivalent forms of TD(γ) updates

$$p_t^s + \gamma (r_t^s + \sum_{q=1}^N p_t^q \psi^q - \sum_{q=1}^N p_{t-1}^q \psi^q) \psi^s; a_t^s$$

Term	Left side	Left side coe .	Right side	Right side coe .
r_t^s	$p_t^s \psi^s; a_t^s$	p_t^s		

Forward view for TD(γ)

Equivalent forms of TD(γ) updates

$$r_t + \gamma V_t - V_{t-1} = \sum_{j=0}^{\infty} \gamma^j (r_{t+j} - V_{t+j-1})$$

$$r_t + \gamma V_t - V_{t-1} = \sum_{j=0}^{\infty} \gamma^j (r_{t+j} - V_{t-1}) - \sum_{j=1}^{\infty} \gamma^j (V_{t+j-1} - V_{t-1})$$

Term	Left side	Left side coe .	Right side	Right side coe .
$r_t + \gamma V_t - V_{t-1}$	$\sum_{j=0}^{\infty} \gamma^j (r_{t+j} - V_{t-1}) - \sum_{j=1}^{\infty} \gamma^j (V_{t+j-1} - V_{t-1})$	$r_t + \gamma V_t - V_{t-1}$	$\sum_{j=0}^{\infty} \gamma^j (r_{t+j} - V_{t-1})$ for $N \geq t$	

Forward view for TD(γ)

Equivalent forms of TD(γ) updates

$$r_t + \gamma V_{\theta}(s_t) - V_{\theta}(s_{t-1}) = \sum_{j=1}^N \phi_j(s_t) \left[p_j - \sum_{i=1}^N \phi_i(s_{t-1}) w_i \right]$$

Term	Left side	Left side coe .	Right side	Right side coe .
$r_t + \gamma V_{\theta}(s_t)$	$r_t + \gamma \sum_{j=1}^N \phi_j(s_t) w_j$	p_j	$\sum_{i=1}^N \phi_i(s_{t-1}) w_i$ for N i t	$p_1 \quad \phi_1(s_{t-1}) \quad \phi_2(s_{t-1}) \quad \dots \quad \phi_N(s_{t-1})$
$V_{\theta}(s_t)$				
$(t \geq 0)$				

Forward view for TD(γ)

Equivalent forms of TD(γ) updates

$$V_{p,q}^{t+1} = \gamma V_{p,q}^t + \delta (r_{p,q}^t + V_{p,q}^t - V_{p,q}^{t-1})$$

Term	Left side	Left side coe .	Right side	Right side coe .
$r_{p,q}^{t+1}$	$\delta (r_{p,q}^t + V_{p,q}^t - V_{p,q}^{t-1})$	$\delta (r_{p,q}^t - V_{p,q}^{t-1})$	$\gamma V_{p,q}^t$ for $N_j = t$	$\gamma (V_{p,q}^{t-1} - V_{p,q}^{t-2} + \dots - V_{p,q}^0)$
$V_{p,q}^{t+1}$ ($t \geq 0$)	$\gamma V_{p,q}^t + \delta (r_{p,q}^t + V_{p,q}^t - V_{p,q}^{t-1})$	$\gamma V_{p,q}^t + \delta (r_{p,q}^t - V_{p,q}^{t-1})$		

Forward view for TD(γ)

Equivalent forms of TD(γ) updates

$$V_{p,q}^{t+1} = \gamma V_{p,q}^t + \delta (r_{p,q}^t + V_{p,q}^t - V_{p,q}^{t-1})$$

Term	Left side	Left side coe .	Right side	Right side coe .
$r_{p,q}^{t+1}$	$\delta (r_{p,q}^t + V_{p,q}^t - V_{p,q}^{t-1})$	$\delta (r_{p,q}^t - V_{p,q}^{t-1})$	$\delta (r_{p,q}^t + V_{p,q}^t)$ for $N_j = t$	$\delta (r_{p,q}^t - V_{p,q}^{t-1})$
$V_{p,q}^{t+1}$ ($t \geq 0$)	$\gamma V_{p,q}^t + \delta (r_{p,q}^t + V_{p,q}^t - V_{p,q}^{t-1})$	$\gamma V_{p,q}^t - V_{p,q}^{t-1}$		

Forward view for TD(γ)

Equivalent forms of TD(γ) updates

$$V_{p,q}^{t+1} = \gamma V_{p,q}^t + (1-\gamma) [r_{p,q}^t + \sum_{q'} p_{N,q'}^t V_{p,q'}^t - \sum_{q'} p_{N,q'}^t V_{p,q}^t]$$

Term	Left side	Left side coe .	Right side	Right side coe .
$r_{p,q}^t$	$p_{1,q} V_{p,q}^{t+1}$	$p_{1,q}$	$p_{N,q} [r_{p,q}^t + \sum_{q'} p_{N,q'}^t V_{p,q'}^t - \sum_{q'} p_{N,q'}^t V_{p,q}^t]$ for $N_j = t$	$p_{1,q} - \sum_{q'} p_{N,q'}^t p_{1,q}$
$V_{p,q}^{t+1}$ ($t \geq 0$)	$p_{1,q} V_{p,q}^{t+1} + \sum_{q'} p_{N,q'}^t V_{p,q'}^t - \sum_{q'} p_{N,q'}^t V_{p,q}^t$	$p_{1,q} - \sum_{q'} p_{N,q'}^t p_{1,q}$	$\sum_{q'} p_{N,q'}^t V_{p,q'}^t - \sum_{q'} p_{N,q'}^t V_{p,q}^t$	$\sum_{q'} p_{N,q'}^t p_{1,q} - \sum_{q'} p_{N,q'}^t p_{1,q}$

Forward view for TD()

Equivalent forms of TD() updates

$$r_{ptq} + \sum_{s=0}^{\infty} \gamma^s p_{1q} \psi^s; a^{ptq} - p_{1q} + \sum_{s=0}^{\infty} \gamma^s \sum_{N=1}^{N-1} p_{Nq} \psi^s; aq$$

Term	Left side	Left side coefficients	Right side	Right side coefficients
r_{ptq}	$p_{1q} \psi^s; a^{ptq}$	$p_{1q} \gamma^s$	$p_{Nq} \psi^s; aq$ for $N = 1, \dots, t$	$p_{1q} \gamma^s + \sum_{N=1}^{N-1} p_{Nq} \gamma^s$
$V_{ps^{ptq}}(t_i = 0)$	$p_{1q} \psi^{s-1}; a^{pt-1q}$ $p_{1q} \psi^s; a^{ptq}$	$p_{1q} \gamma^{s-1} + p_{1q} \gamma^s$	$p_{tq} \psi^s; aq$	$p_{1q} \gamma^{s-1} + \sum_{N=1}^{N-1} p_{Nq} \gamma^s$
V_{psq}				

Forward view for TD()

Equivalent forms of TD() updates

$$r_{ptq} + \sum_{t=0}^{\infty} \gamma^t p_{1q} \psi^{ptq}; a^{ptq} q - p_{1q} - \sum_{N=1}^{\infty} \gamma^{N-1} p_{Nq} \psi; aq$$

Term	Left side	Left side coef.	Right side	Right side coef.
r_{ptq}	$p_{1q} \psi^{ptq}; a^{ptq} q$	$p_{1q} \gamma^t$	$p_{Nq} \psi; aq$ for $N = 1 \dots t$	$\gamma^{t-1} p_{1q} \gamma^{t-2} \dots \gamma^0 p_{Nq}$
$V_{ps^{ptq}q}$ ($t \geq 0$)	$p_{1q} \psi^{pt-1q}; a^{pt-1q} q$ $p_{1q} \psi^{ptq}; a^{ptq} q$	$p_{1q} \gamma^{t-1} - p_{1q} \gamma^t$	$p_{tq} \psi; aq$	$p_{1q} \gamma^{t-1} p_{tq}$
V_{psq}	$p_{1q} \psi; aq$	1		

Forward view for TD()

Equivalent forms of TD() updates

$$r_{ptq} + \sum_{t=0}^{\infty} \gamma^t p_{1q} \psi^{ptq}; a^{ptq} q - p_{1q} - \sum_{N=1}^{\infty} \gamma^{N-1} p_{Nq} \psi; aq$$

Term	Left side	Left side coef.	Right side	Right side coef.
r_{ptq}	$p_{1q} \psi^{ptq}; a^{ptq} q$	$p_{1q} q^t$	$p_{Nq} \psi; aq$ for $N_j = t$	$t p_{1q} \quad q p_{t-1} \quad t-1 \quad t-2 \quad \dots : q$
$V_{\psi^{ptq} q}$ ($t_j = 0$)	$p_{1q} \psi^{pt-1q}; a^{pt-1q} q$ $p_{1q} \psi^{ptq}; a^{ptq} q$	$p_{1q} q^{t-1} \quad p_{1q} q^t$	$p_{tq} \psi; aq$	$p_{1q} \quad q p_{t-1} \quad t-1 \quad t-2 \quad \dots : q$
$V_{\psi q}$	$p_{1q} \psi; aq$	1	all	$p_{1q} \quad q p_{1q} \quad 2 \quad \dots : q$

Generalized Advantage Estimation (GAE)

$$GAE_{p_S; a_Q} : \sum_{t=0}^T p_{1Q} p_S^{ptQ}; a^{ptQ} q$$

(trace is zeroed when future is not available)

Generalized Advantage Estimation (GAE)

$$GAE_{p_S; a_Q} = \sum_{t=0}^T \gamma^t p_{1-Q} p_S^{pt_Q}; a^{pt_Q} Q$$

(trace is zeroed when future is not available)

How to compute in practice:

$$GAE_{p_S; a_Q} = p_{1-Q} p_S; a_Q - p_{1-Q} \text{done}_{t-1} Q + \gamma GAE_{p_S; a_Q}$$

Proximal Policy Optimization: implementation matters

Key elements:

- × Clipped policy loss
- × Clipped critic loss
- × GAE

Pipeline details:

- ! Advantage normalization in mini-batches
- No KL regularization
- Entropy loss

¹divided by running std of collected cumulative rewards

²can be critical in continuous control

Other hacks:

- ! Reward normalization¹ and clipping
- Observations normalization and clipping²
- Orthogonal initialization of layers
- (clipping parameter) annealing

Standard tricks:

- Adam, learning rate annealing
- Tanh activation functions
- ! Gradient clipping

Full Pipeline: pt.1

Proximal Policy Optimization (PPO)

Initialize $\pi, \mu, \sigma, \gamma, \beta, \epsilon, \tau$

Full Pipeline: pt.1

Proximal Policy Optimization (PPO)

Initialize π, μ, σ, V

for $k = 0; 1; 2; \dots$:

- collect several rollouts $s_0; a_0; r_0; s_1; done_1; a_1; \dots; s_N; done_N$ using π, μ, σ
store probabilities of selected actions as $\pi^{old} \pi(a_t | s_t; \mu, \sigma)$
store critic output as $V^{old} V(s_t; \mu, \sigma)$

Full Pipeline: pt.1

Proximal Policy Optimization (PPO)

Initialize $\pi(a|s); q; V; \psi; \sigma$

for $k = 0; 1; 2; \dots$:

- collect several rollouts $s_0; a_0; r_0; s_1; \dots; s_N; a_N$ using $\pi(a|s); q$
store probabilities of selected actions as $\pi^{\text{old}}(a_t|s_t); \pi(a_t|s_t); q$
store critic output as $V^{\text{old}}(s_t); V(s_t); q$
- compute 1-step errors: $\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t); q$

Full Pipeline: pt.1

Proximal Policy Optimization (PPO)

Initialize $\pi, \mu, \sigma, V, \gamma$

for $k = 0; 1; 2; \dots$:

- collect several rollouts $s_0; a_0; r_0; s_1; done_1; a_1; \dots; s_N; done_N$ using π, μ, σ
store probabilities of selected actions as $\pi^{old}(a_t | s_t)$
store critic output as $V^{old}(s_t)$
- compute 1-step errors: $\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$
- compute GAE advantage estimations: $A_t = \sum_{l=0}^{\infty} \gamma^l \delta_{t+l}$
- for t from $N-2$ to 0:
 - $A_t = \delta_t + \gamma V(s_{t+1}) - V(s_t)$

Full Pipeline: pt.1

Proximal Policy Optimization (PPO)

Initialize $\pi, \mu, \sigma, V, \gamma$

for $k = 0; 1; 2; \dots$:

- collect several rollouts $s_0; a_0; r_0; s_1; done_1; a_1; \dots; s_N; done_N$ using π, μ, σ
store probabilities of selected actions as $\pi^{old}(a_t | s_t)$
store critic output as $V^{old}(s_t)$
- compute 1-step errors: $\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$
- compute GAE advantage estimations: $A_t = \sum_{l=0}^{\infty} \gamma^l \delta_{t+l}$
- for t from $N-2$ to 0:
 - $A_t = \delta_t + \gamma V(s_{t+1}) - V(s_t) + \gamma A_{t+1}$

Full Pipeline: pt.1

Proximal Policy Optimization (PPO)

Initialize $\pi, \mu, \sigma, V; \theta, \phi$

for $k = 0; 1; 2; \dots$:

- collect several rollouts $s_0; a_0; r_0; s_1; done_1; a_1; \dots; s_N; done_N$ using $\pi, \mu, \sigma; \theta, \phi$
store probabilities of selected actions as $\pi^{old} p(a_t | s_t; \theta, \phi)$
store critic output as $V^{old} p(s_t; \theta, \phi)$
- compute 1-step errors: $\delta_t = r_t + \gamma V(s_{t+1}; \theta, \phi) - V(s_t; \theta, \phi)$
- compute GAE advantage estimations: $A_t = \sum_{l=0}^{\infty} \gamma^l \delta_{t+l}$
- for t from $N-2$ to 0:
 - $\hat{A}_t = A_t + \lambda (V(s_t; \theta, \phi) - V(s_{t+1}; \theta, \phi))$
- compute critic targets: $y_t = V(s_t; \theta, \phi) + \lambda A_t$

Full Pipeline: pt.1

Proximal Policy Optimization (PPO)

Initialize $\pi, \mu, \sigma, V, \gamma$

for $k = 0; 1; 2; \dots$:

- collect several rollouts $s_0; a_0; r_0; s_1; done_1; a_1; \dots; s_N; done_N$ using π, μ, σ
store probabilities of selected actions as $\pi^{old}(a_t | s_t)$
store critic output as $V^{old}(s_t)$
- compute 1-step errors: $\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$
- compute GAE advantage estimations: $A_t = \sum_{l=0}^{\infty} \gamma^l \delta_{t+l}$
- for t from $N-2$ to 0:
 - $A_t = \delta_t + \gamma V(s_{t+1}) - V(s_t) + \gamma A_{t+1}$
- compute critic targets: $y_t = V(s_t) + A_t$
- construct dataset of $(s_t, a_t, A_t, y_t, \pi^{old}(a_t | s_t), V^{old}(s_t))$

Proximal Policy Optimization (PPO) -- cont.

- go through dataset n_epochs times, sampling mini-batches of size B ; for each mini-batch:

Proximal Policy Optimization (PPO) -- cont.

- go through dataset n_epochs times, sampling mini-batches of size B ; for each mini-batch:
 - normalize $GAE_{ps;aq}$ in the batch by subtracting mean and dividing by std

Full Pipeline: pt.II

Proximal Policy Optimization (PPO) -- cont.

- go through dataset n_epochs times, sampling mini-batches of size B ; for each mini-batch:
 - normalize $GAE_{p_{s;a;q}}$ in the batch by subtracting mean and dividing by std
 - compute importance sampling weights:

$$p_{s;a;q} : \frac{p_{a|s;q}}{old_{p_{a|s;q}}}, \quad clip_{p_{s;a;q}} = clip(p_{s;a;q}, 1-\epsilon, 1+\epsilon)$$

Full Pipeline: pt.II

Proximal Policy Optimization (PPO) -- cont.

- go through dataset n_epochs times, sampling mini-batches of size B ; for each mini-batch:
 - normalize $GAE_{ps;a;q}$ in the batch by subtracting mean and dividing by std
 - compute importance sampling weights:

$$w_{ps;a;q} = \frac{p_{\pi}(a|s;q)}{p_{\pi}^{old}(a|s;q)}, \quad \text{clip}_{\epsilon}(w_{ps;a;q}) = \min\left(\frac{w_{ps;a;q}}{1+\epsilon}, \frac{w_{ps;a;q}}{1-\epsilon}\right)$$

- update actor:

$$L_1(\pi; \theta) = \sum_{s,a} w_{ps;a;q} GAE_{ps;a;q} \log \frac{p_{\pi}(a|s;q)}{p_{\pi}^{old}(a|s;q)}$$

$$\mathbb{E}_{s,a} \left[r + \frac{1}{B} \sum_{s,a} \min_p L_1(\pi; \theta); L_2(\pi; \theta) \right]$$

Full Pipeline: pt.II

Proximal Policy Optimization (PPO) -- cont.

- go through dataset n_epochs times, sampling mini-batches of size B ; for each mini-batch:
 - normalize $GAE_{ps;a;q}$ in the batch by subtracting mean and dividing by std
 - compute importance sampling weights:

$$p_{s;a;q} : \frac{p_{old}|s;q}{p_{new}|s;q}, \quad \text{clip}_{ps;a;q} = \min\left(\frac{p_{new}|s;q}{p_{old}|s;q}, 1\right), \max\left(\frac{p_{new}|s;q}{p_{old}|s;q}, 1\right)$$

- update actor:

$$L_1_{ps;a;q} : p_{ps;a;q} GAE_{ps;a;q}, \quad L_2_{ps;a;q} : \text{clip}_{ps;a;q} GAE_{ps;a;q}$$

$$\mathbb{E}_s \left[r \frac{1}{B} \sum_{s;a} \min_p L_1_{ps;a;q}; L_2_{ps;a;q} \right]$$

- update critic:

$$\text{Loss}_1_{ps;q} : y_{psq} - V_{psq}^2$$

$$\text{Loss}_2_{ps;q} : y_{psq} - V_{psq}^{\text{old}} - \text{clip}(V_{psq} - V_{psq}^{\text{old}}, \epsilon, \epsilon)$$

$$\mathbb{E}_s \left[r \frac{1}{B} \sum_s \max_p \text{Loss}_1_{ps;q}; \text{Loss}_2_{ps;q} \right]$$

