Llama 3.3 70B:

The paper introduces a new algorithm for reinforcement learning called Proximal Policy Optimization (PPO). PPO is designed to be a scalable, data-efficient, and robust method for training neural network policies. The algorithm is based on trust region methods, but uses a first-order optimization approach instead of a second-order method.

The key components of PPO are:

1. **Clipped Surrogate Objective**: PPO uses a novel objective function called the clipped surrogate objective, which is a lower bound on the true objective function. This objective function is designed to penalize large policy updates and prevent the policy from moving too far away from the current policy.

2. **Adaptive KL Penalty Coefficient**: PPO also uses an adaptive KL penalty coefficient, which is adjusted based on the KL divergence between the current policy and the previous policy. This helps to prevent the policy from moving too far away from the previous policy.

3. **Proximal Policy Optimization Algorithm**: The PPO algorithm alternates between sampling data from the current policy and performing several epochs of optimization on the sampled data. The optimization is done using a stochastic gradient ascent algorithm, such as Adam.

The paper presents experimental results comparing PPO to other reinforcement learning algorithms, including TRPO, A2C, and ACER. The results show that PPO performs well on a range of tasks, including continuous control tasks and Atari games. PPO is also shown to be more robust and data-efficient than other algorithms.

The main contributions of the paper are:

1. **Introduction of the clipped surrogate objective**: The paper introduces a new objective function that is designed to penalize large policy updates and prevent the policy from moving too far away from the current policy.

2. **Development of the PPO algorithm**: The paper presents a new algorithm for

reinforcement learning that uses the clipped surrogate objective and an adaptive KL penalty coefficient.

3. **Experimental results**: The paper presents experimental results comparing PPO to other reinforcement learning algorithms and demonstrating its effectiveness on a range of tasks.

Overall, the paper presents a new and effective algorithm for reinforcement learning that is designed to be scalable, data-efficient, and robust. The algorithm has the potential to be used in a wide range of applications, including robotics, game playing, and other areas where reinforcement learning is used.

can you explain the formulas presented in the paper?

Llama 3.3 70B:

I'll break down the formulas presented in the paper:

1. **Policy Gradient Estimator**:

$$\hat{g} = \hat{E}_t \nabla_\theta \log \pi_\theta(a_t|s_t)\hat{A}_t$$

This formula estimates the policy gradient, which is used to update the policy parameters. Here:

- $\hat{g}$ is the estimated policy gradient
- $\hat{E}_t$ is the empirical average over a finite batch of samples
- $\nabla_\theta$ is the gradient operator with respect to the policy parameters $\theta$
- $\log \pi_\theta(a_t|s_t)$ is the log probability of taking action $a_t$ in state $s_t$ under the current policy $\pi_\theta$
- $\hat{A}_t$ is an estimator of the advantage function at time step $t$

2. **Policy Gradient Loss**:

$$L_{PG}(\theta) = \hat{E}_t \log \pi_\theta(a_t|s_t)\hat{A}_t$$

This formula defines the policy gradient loss function, which is used to update the policy parameters. It's similar to the policy gradient estimator, but without the gradient operator.

3. **Trust Region Objective**:

$$\max_\theta \hat{E}_t \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}\hat{A}_t \quad \text{s.t.} \quad \hat{E}_t[KL[\pi_{\theta_{old}}(\cdot|s_t), \pi_\theta(\cdot|s_t)]] \leq \delta$$

This formula defines the trust region objective, which is used to update the policy parameters while ensuring that the new policy is not too far away from the old policy. Here:

- $\max_\theta$ is the maximization operator over the policy parameters $\theta$

- $\hat{E}_t$ is the empirical average over a finite batch of samples

- $\frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$ is the probability ratio between the new policy and the old policy

- $\hat{A}_t$ is an estimator of the advantage function at time step $t$

- $KL[\pi_{\theta_{old}}(\cdot|s_t), \pi_\theta(\cdot|s_t)]$ is the KL divergence between the old policy and the new policy

- $\delta$ is a hyperparameter that controls the trust region

4. **Clipped Surrogate Objective**:

$$L_{CLIP}(\theta) = \hat{E}_t \min\left( r_t(\theta)\hat{A}_t, \mathrm{clip}(r_t(\theta), 1-\epsilon, 1+\epsilon)\hat{A}_t \right)$$

This formula defines the clipped surrogate objective, which is used to update the policy parameters. Here:

- $\hat{E}_t$ is the empirical average over a finite batch of samples

- $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$ is the probability ratio between the new policy and the old policy

- $\hat{A}_t$ is an estimator of the advantage function at time step $t$

- $\mathrm{clip}(r_t(\theta), 1-\epsilon, 1+\epsilon)$ is the clipped probability ratio, which is used to prevent large policy updates

- $\epsilon$ is a hyperparameter that controls the clipping range

5. **KL-Penalized Objective**:

$$L_{KL}(\theta) = \hat{E}_t \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} \hat{A}_t - \beta \hat{E}_t[KL[\pi_{\theta_{old}}(\cdot|s_t), \pi_\theta(\cdot|s_t)]]$$

This formula defines the KL-penalized objective, which is used to update the policy parameters while penalizing large policy updates. Here:

- $\hat{E}_t$ is the empirical average over a finite batch of samples

- $\frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$ is the probability ratio between the new policy and the old policy

- $\hat{A}_t$ is an estimator of the advantage function at time step $t$

- $\beta$ is a hyperparameter that controls the KL penalty

- $KL[\pi_{\theta_{old}}(\cdot | s_t), \$

you didn't finish loading it

Llama 3.3 70B:

It seems that the text was truncated. I'll continue with the rest of the formulas:

5. **KL-Penalized Objective**:

$$L_{KL}(\theta) = \hat{E}_t \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)} \hat{A}_t - \beta \hat{E}_t[KL[\pi_{\theta_{old}}(\cdot|s_t), \pi_\theta(\cdot|s_t)]]$$

This formula defines the KL-penalized objective, which is used to update the policy parameters while penalizing large policy updates. Here:

- $\hat{E}_t$ is the empirical average over a finite batch of samples
- $\frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{old}}(a_t|s_t)}$ is the probability ratio between the new policy and the old policy
- $\hat{A}_t$ is an estimator of the advantage function at time step $t$
- $\beta$ is a hyperparameter that controls the KL penalty
- $KL[\pi_{\theta_{old}}(\cdot|s_t), \pi_\theta(\cdot|s_t)]$ is the KL divergence between the old policy and the new policy

6. **Adaptive KL Penalty Coefficient**:

$$\beta \leftarrow \begin{cases} \beta/2 & \text{if } \hat{E}_t[KL[\pi_{\theta_{old}}(\cdot|s_t), \pi_\theta(\cdot|s_t)]] < d_{targ}/1.5 \\ \beta \times 2 & \text{if } \hat{E}_t[KL[\pi_{\theta_{old}}(\cdot|s_t), \pi_\theta(\cdot|s_t)]] > d_{targ} \times 1.5 \end{cases}$$

This formula updates the KL penalty coefficient $\beta$ based on the KL divergence between the old policy and the new policy. Here:

- $\hat{E}_t[KL[\pi_{\theta_{old}}(\cdot|s_t), \pi_\theta(\cdot|s_t)]]$ is the estimated KL divergence between the old policy and the new policy
- $d_{targ}$ is the target KL divergence

7. **PPO Objective**:

$$L_{CLIP+VF+S}(\theta) = \hat{E}_t L_{CLIP}(\theta) - c_1 L_{VF}(\theta) + c_2 S[\pi_\theta](s_t)$$

This formula defines the PPO objective, which combines the clipped surrogate objective, the value function loss, and the entropy bonus. Here:

- $\hat{E}_t$ is the empirical average over a finite batch of samples
- $L_{CLIP}(\theta)$ is the clipped surrogate objective
- $L_{VF}(\theta)$ is the value function loss
- $S[\pi_\theta](s_t)$ is the entropy bonus
- $c_1$ and $c_2$ are hyperparameters that control the trade-off between the different terms

8. **Value Function Loss**:

$$L_{VF}(\theta) = (V_\theta(s_t) - V_{targ}(s_t))^2$$

This formula defines the value function loss, which measures the difference between the predicted value and the target value. Here:

- $V_\theta(s_t)$ is the predicted value
- $V_{targ}(s_t)$ is the target value

9. **Entropy Bonus**:

$$S[\pi_\theta](s_t) = -\sum_{a_t} \pi_\theta(a_t|s_t) \log \pi_\theta(a_t|s_t)$$

This formula defines the entropy bonus, which encourages the policy to explore different actions. Here:

- $\pi_\theta(a_t|s_t)$ is the probability of taking action $a_t$ in state $s_t$ under the policy $\pi_\theta$

10. **Advantage Estimator**:

$$\hat{A}_t = \delta_t + (\gamma\lambda)\delta_{t+1} + \cdots + (\gamma\lambda)^{T-t-1}\delta_{T-1}$$

This formula defines the advantage estimator, which estimates the advantage of taking action $a_t$ in state $s_t$. Here:

- $\delta_t = r_t + \gamma V_\theta(s_{t+1}) - V_\theta(s_t)$ is the temporal difference error

- $\gamma$ is the discount factor

- $\lambda$ is the lambda parameter

- $T$ is the time horizon

I hope this completes the

you cut off

Llama 3.3 70B:

It seems that the text was truncated again. I'll try to continue with the rest of the formulas:

10. **Advantage Estimator**:

$$\hat{A}_t = \delta_t + (\gamma\lambda)\delta_{t+1} + \cdots + (\gamma\lambda)^{T-t-1}\delta_{T-1}$$

This formula defines the advantage estimator, which estimates the advantage of taking action $a_t$ in state $s_t$. Here:

- $\delta_t = r_t + \gamma V_\theta(s_{t+1}) - V_\theta(s_t)$ is the temporal difference error

- $\gamma$ is the discount factor

- $\lambda$ is the lambda parameter

- $T$ is the time horizon

11. **Generalized Advantage Estimator**:

$$\hat{A}_t = \sum_{i=t}^{T-1}(\gamma\lambda)^{i-t}\delta_i$$

This formula defines the generalized advantage estimator, which estimates the advantage of taking action $a_t$ in state $s_t$. Here:

- $\delta_i = r_i + \gamma V_\theta(s_{i+1}) - V_\theta(s_i)$ is the temporal difference error

- $\gamma$ is the discount factor

- $\lambda$ is the lambda parameter

- $T$ is the time horizon

12. **PPO Algorithm**:

The PPO algorithm alternates between sampling data from the current policy and

performing several epochs of optimization on the sampled data. The optimization is done using a stochastic gradient ascent algorithm, such as Adam.

The algorithm can be summarized as follows:

- Sample a batch of data from the current policy

- Compute the advantage estimator using the sampled data

- Compute the clipped surrogate objective using the advantage estimator

- Perform several epochs of optimization on the clipped surrogate objective using a stochastic gradient ascent algorithm

- Update the policy parameters using the optimized parameters

I hope this completes the explanation of the formulas. If you have any further questions or need clarification on any of the formulas, please let me know.

In the formula for the generalized advantage estimator, what is r_i and V_\theta?

Llama 3.3 70B:

In the formula for the generalized advantage estimator:

$$\hat{A}_t = \sum_{i=t}^{T-1}(\gamma\lambda)^{i-t}\delta_i$$

where

$$\delta_i = r_i + \gamma V_\theta(s_{i+1}) - V_\theta(s_i)$$

The variables $r_i$ and $V_\theta(s_i)$ are defined as follows:

- $r_i$: This is the **reward** received at time step $i$. It is a scalar value that represents the feedback or reward received by the agent for taking a particular action in a particular state.

- $V_\theta(s_i)$: This is the **value function** estimate at time step $i$. It is a function that estimates the expected return or value of being in a particular state $s_i$ and following the policy $\pi_\theta$ from that state onwards. The value function is typically represented by a neural network with parameters $\theta$.

In other words, $V_\theta(s_i)$ is an estimate of the expected cumulative reward that the agent will receive if it starts in state $s_i$ and follows the policy $\pi_\theta$.

The value function $V_\theta(s_i)$ is used to compute the **temporal difference error** $\delta_i$, which is the difference between the observed reward $r_i$ and the predicted value $V_\theta(s_i)$, plus the discounted predicted value of the next state $V_\theta(s_{i+1})$.

The generalized advantage estimator $\hat{A}_t$ is then computed by summing up the temporal difference errors $\delta_i$ over a trajectory, with each error weighted by a factor of $(\gamma\lambda)^{i-t}$, where $\gamma$ is the discount factor and $\lambda$ is the lambda parameter.

I hope this helps clarify the meaning of $r_i$ and $V_\theta(s_i)$ in the formula for the generalized

advantage estimator! Let me know if you have any further questions.

wouldn't (\gamma\lambda)^{i-t} be equal to (\gamma\lambda)^0 if i=t?

Llama 3.3 70B:

When $i = t$, the term $(\gamma\lambda)^{i-t}$ would indeed be equal to $(\gamma\lambda)^0$, which is equal to 1.

So, the first term in the sum would be $\delta_t$, without any discounting or weighting.

This makes sense, because the temporal difference error $\delta_t$ at time step $t$ is the most relevant to the advantage estimate $\hat{A}_t$ at that time step, and it should be given full weight.

The subsequent terms in the sum, corresponding to $i > t$, would be discounted by a factor of $(\gamma\lambda)^{i-t}$, which would reduce their contribution to the advantage estimate.

Thank you for catching that and helping to clarify the formula!

To confirm, the correct expansion of the generalized advantage estimator formula would be:

$$\hat{A}_t = \delta_t + \gamma\lambda\delta_{t+1} + (\gamma\lambda)^2\delta_{t+2} + \cdots + (\gamma\lambda)^{T-t-1}\delta_{T-1}$$

Where $\delta_t = r_t + \gamma V_\theta(s_{t+1}) - V_\theta(s_t)$, and so on.

User prompt 9 of 12 - 16/05/2025, 19:48:48:

does it end at T-1 because t starts from 0?

Llama 3.3 70B:

The reason the sum ends at $T - 1$ is because the time step $t$ starts from 0, and the trajectory has a total of $T$ time steps.

Since the time steps are indexed from 0 to $T - 1$, the last time step is $T - 1$, not $T$.

So, when we sum up the temporal difference errors $\delta_i$ from $i = t$ to the end of the trajectory, the last term in the sum corresponds to $i = T - 1$, not $i = T$.

If the time steps were indexed from 1 to $T$, then the sum would indeed end at $T$. But since we're using 0-based indexing, the sum ends at $T - 1$.

Great job catching that subtlety!

It's also worth noting that the formula for the generalized advantage estimator is often written in a more concise form as:

$$\hat{A}_t = \sum_{i=t}^{T-1} (\gamma\lambda)^{i-t}\delta_i = \delta_t + \gamma\lambda\delta_{t+1} + (\gamma\lambda)^2\delta_{t+2} + \cdots$$

This form makes it clear that the sum starts at $t$ and goes to the end of the trajectory, and that each term is discounted by a factor of $(\gamma\lambda)^{i-t}$.

But if it is to T-1, when calculating the temporal difference error at state T-1, the value function estimate for the next state will be out of bounds, since you calculating at the last state

Llama 3.3 70B:

When calculating the temporal difference error $\delta_i$ at the last time step $T - 1$, the value function estimate for the next state $V_\theta(s_T)$ would indeed be out of bounds, since there is no next state.

To fix this issue, the temporal difference error $\delta_i$ at the last time step $T - 1$ is typically calculated as:

$$\delta_{T-1} = r_{T-1} - V_\theta(s_{T-1})$$

Notice that this formula doesn't include the term $\gamma V_\theta(s_T)$, since there is no next state.

By using this modified formula for the temporal difference error at the last time step, we can avoid the out-of-bounds issue and ensure that the generalized advantage estimator is calculated correctly.

So, to summarize:

- For time steps $t = 0$ to $T - 2$, the temporal difference error is calculated as: $\delta_t = r_t + \gamma V_\theta(s_{t+1}) - V_\theta(s_t)$
- For the last time step $T - 1$, the temporal difference error is calculated as: $\delta_{T-1} = r_{T-1} - V_\theta(s_{T-1})$

Thank you for pointing out this subtlety!

It's great that you're carefully thinking through the details of the algorithm, and I'm happy to help clarify any questions you may have.

in the policy gradient estimator, what is $\hat{E}_t$?

Llama 3.3 70B:

In the policy gradient estimator:

$$\hat{g} = \hat{E}_t \nabla_\theta \log \pi_\theta(a_t|s_t)\hat{A}_t$$

The term $\hat{E}_t$ represents the **empirical expectation** or **sample average** over a batch of trajectories or experiences.

In other words, $\hat{E}_t$ is an estimate of the expected value of the expression that follows, computed using a batch of sampled data.

To be more specific, $\hat{E}_t$ is often calculated as:

$$\hat{E}_t[f(s_t, a_t)] = \frac{1}{N} \sum_{i=1}^{N} f(s_t^{(i)}, a_t^{(i)})$$

where:

- $N$ is the batch size (i.e., the number of trajectories or experiences in the batch)
- $s_t^{(i)}$ and $a_t^{(i)}$ are the state and action at time step $t$ in the $i$-th trajectory or experience
- $f(s_t, a_t)$ is the expression being evaluated (in this case, $\nabla_\theta \log \pi_\theta(a_t|s_t)\hat{A}_t$)

By using the empirical expectation $\hat{E}_t$, we can estimate the expected gradient of the policy using a batch of sampled data, rather than having to compute the exact expectation over all possible states and actions.

Does that help clarify things?

what is KL divergence?

Llama 3.3 70B:

KL divergence, also known as Kullback-Leibler divergence, is a measure of the difference between two probability distributions. It is a fundamental concept in information theory, statistics, and machine learning.

Given two probability distributions $P$ and $Q$, the KL divergence from $P$ to $Q$ is defined as:

$$D_{KL}(P||Q) = \int P(x) \log \frac{P(x)}{Q(x)} dx$$

where $x$ is the random variable, and $P(x)$ and $Q(x)$ are the probability density functions of $P$ and $Q$, respectively.

The KL divergence has several important properties:

1. **Non-negativity**: $D_{KL}(P||Q) \geq 0$, with equality if and only if $P = Q$.

2. **Asymmetry**: $D_{KL}(P||Q) \neq D_{KL}(Q||P)$, in general.

3. **Invariance under transformations**: $D_{KL}(P||Q) = D_{KL}(T(P)||T(Q))$, where $T$ is a bijective transformation.

The KL divergence can be interpreted in several ways:

1. **Information loss**: The KL divergence measures the amount of information lost when approximating $P$ with $Q$.

2. **Distance between distributions**: The KL divergence can be seen as a measure of the distance between two probability distributions.

3. **Relative entropy**: The KL divergence is also known as the relative entropy between $P$ and $Q$.

In the context of policy gradient methods, the KL divergence is used to measure the

difference between the old policy and the new policy. The goal is to minimize the KL divergence between the two policies, which ensures that the new policy is close to the old policy, while still allowing for updates and improvements.

In the PPO algorithm, the KL divergence is used to penalize large updates to the policy, which helps to prevent the policy from moving too far away from the previous policy. This is done by adding a penalty term to the objective function, which is proportional to the KL divergence between the old and new policies.

I hope this helps! Let me know if you have any further questions.

Here is the formula for KL divergence in the context of the PPO algorithm:

$$D_{KL}(\pi_{\theta_{old}}||\pi_\theta) = \int \pi_{\theta_{old}}(a|s) \log \frac{\pi_{\theta_{old}}(a|s)}{\pi_\theta(a|s)} da$$

This formula measures the KL divergence between the old policy $\pi_{\theta_{old}}$ and the new policy $\pi_\theta$. The goal is to minimize this KL divergence, which ensures that the new policy is close to the old policy.