

# Policy gradient based methods: concepts

Recall the Policy Gradient Theorem

- formulated with single reward  $\text{rewseq}(\tau)$  for the entire trajectory

$$\nabla_{\theta} J(\theta) = \text{Exp} \tau \sim \text{pr} \theta \sum_{=0}^{|\tau|} \nabla_{\theta} \log \pi(\text{actseq}_{\tau}, | \text{stateseq}_{\tau}) \text{rewseq}(\tau)$$

- formulated with intermediate rewards

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_{=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(\text{actseq}_{\tau}, | \text{stateseq}_{\tau,t}) G_{\tau} \right]$$

Policy based methods update parameters  $\theta$

- in the direction (gradient)
- that increases expected return

## Advantage vs Return/Reward: Baselines

Suppose all returns/rewards are positive.

The Policy Gradient Theorem would suggest that we update parameters

- to favor *all* actions
- with positive gradients

But consider two actions at some step

- one with a *very large* positive return/reward
- one with a *very small* positive return/reward

We probably want to favor the action with large reward.

Thus, we want to *relativize* the return/reward by comparing it to a *baseline* for the state.

The relativized value of an action is called the *advantage*.

Often, the advantage is derived from the action's return/reward by

- subtracting a *baseline*
- where the baseline value
  - is a function of *all* possible actions from the state
  - usually  $\text{\textbackslash statevalfun}_\pi(\text{\textbackslash stateseq})$ , the value of the state as would be computed by a Value function

For example: consider a baseline that is the average return/reward from the state

- above average returns have positive advantage
- below average returns have a negative advantage

Thus, in the case of all returns/rewards being positive

- we favor actions that result in positive advantages

Subtracting a baseline from the return/reward has other advantages

- reduces the *magnitude* of the parameter update
  - smoother training
- reduces "noise"
  - the baseline value for a state is common to *all actions* from the state
  - so the advantage is relative to the signal *only* from the action itself

| Aspect             | Effect of Subtracting Baseline                          |
|--------------------|---|
| Noise cancellation | Removes expected (common) return, reducing fluctuations |
| Bias introduction  | None, baseline does not depend on action                |
| Zero-centering     | Advantages centered around zero, stabilizing updates    |

# Unified Policy Gradient Formulation: Advantage Definitions

Using the concept of Advantage, we can rewrite the Policy Gradient Theorem

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta}(\text{\textbackslash actseq}_{\tau, t} | \text{\textbackslash stateseq}_{\tau, t}) \text{\textbackslash advseq}_{\tau, t} \right]$$

The form of the advantage might vary depending on whether there is

- a single per-trajectory reward
- or intermediate rewards.

We can characterize many policy-based methods

- by defining their advantage calculation  
 $\text{\textbackslash advseq}_{\tau, t}$

We thus refer to the above formulation as the Unified Policy Gradient Formulation.

# Surrogate Loss

Rather than

- maximization of return

Policy methods often switch to

- minimization of a *surrogate loss*

The surrogate loss often **imposes constraints** on the derived policy and on the solution process.

The Policy Gradient Theorem provides

- the *theoretical* basis for an optimal policy
- using the Returns or Advantages

The Surrogate Loss provides

- the *practical* objective for finding the optimal policy
- subject to practical constraints: stability, convergence

The actual implementation of our models will usually revert to the Minimization of Surrogate Loss formulation.

Some practical constraints limit policy changes

- prevents drastic policy shifts after a parameter update
- by constraining the new policy to be close to a "reference" policy

Other constraints try to promote stability and convergence in the solution process

- limiting the magnitude of a gradient update
- promoting low variance of the gradient updates

## Constraints

| Reason                            | Explanation  |
|-----------------------------------|--|
| Trust Region Constraints          | Avoid overly large policy updates that destabilize learning  |
| Regularization                    | KL divergence or clipping adds penalty to maintain stability |
| Stability & Convergence           | Ensures smooth, incremental improvement of the policy        |
| Computational Tractability        | Easier to optimize surrogate than raw return objective       |
| Exploration-Exploitation Tradeoff | Clipping controls how much policy can change per step        |

## Relationship between Policy Gradient Theorem and Surrogate Loss

| Aspect            | Policy Gradient Theorem                          | Surrogate Loss   |
|-------------------|--|--|
| Role              | Theoretical formula for expected return gradient | Practical approximation/objective for policy update                        |
| Gradient          | Directly gives unbiased gradient of $J(\theta)$  | Its gradient approximates policy gradient but includes stabilizing terms   |
| Constraints       | None intrinsic; pure gradient formula            | Includes clipping, penalties to enforce trust region and avoid large steps |
| Use in Algorithms | Foundation for policy gradient methods           | Used in PPO-style algorithms to compute safe gradient steps                |
| Goal              | Maximize expected return $J(\theta)$             | Provide stable, incremental improvement proxy                              |

Here is a preview of the Surrogate Losses of the policy-based methods we will explore.

| Method    | Surrogate Return/Objective Used for Gradient | Credit Assignment             |
|-----------|--|-------------------------------|
| REINFORCE | $G_t$ (return-to-go from $t$ )               | Monte Carlo trajectory return |
| PPO       | Clipped ratio times $A_t$ (advantage)        | Advantage-based, stable       |
| GRPO      | Relative advantages over candidate           |                               |

# Policy gradient based methods: Preview

We will show several common policy-based methods.

- REINFORCE
- PPO
- GRPO (relatively new: 2024)

Details of each follow.

As a preview, we show the definition of the Advantage for each.

| Method    | Intermediate Reward Advantage $\backslash\text{advseq}_t$ ,   | Single Trajectory Reward Advantage $\backslash\text{advseq}_\tau$ |
|-----------|---|---|
| REINFORCE | $\backslash\text{advseq}_t = G_t - b_t$                       | $\backslash\text{advseq} = R(\tau) - b$                           |
| PPO       | $\backslash\text{advseq}_t = \hat{A}_t$ (e.g., GAE estimator) | $\backslash\text{advseq} = R(\tau) - b$ applied per step          |
| GRPO      | Relative advantages per candidate                             | Same per-candidate advantage applied per step                     |

## And a comparison:

| Method                                    | Gradient-Based | Main Objective  | Key Characteristics   | Stability & Sample Efficiency  | Typical Application Domains   |
|---|----------------|---|---|--|---|
| PPO (Proximal Policy Optimization)        | Yes            | Maximize expected reward with clipped surrogate objective     | Uses policy gradients with clipping to limit policy update magnitude, balancing exploration and exploitation                | High stability; more sample efficient than vanilla policy gradients; widely used for continuous and discrete control tasks | Robotics, games, continuous control, benchmark RL tasks             |
| DPO (Direct Preference Optimization)      | Yes            | Directly optimize policy based on preference data             | Uses a preference-based loss to train policy without explicit reward modeling; bypasses traditional RL complexities         | Improved stability by leveraging human preferences; avoids some issues of reward misspecification                          | Alignment of language models with human preferences, NLP-focused RL |
| GRPO (Group Relative Policy Optimization) | Yes            | Optimize policy using group-relative advantage estimates      | Does not require a separate value function; updates policy based on relative advantages within a group of candidate outputs | More memory efficient, stable; effective for large-scale policy optimization with reduced critic reliance                  | Training large language models, large-scale policy optimization     |
| REINFORCE                                 | Yes            | Maximize expected cumulative reward by direct policy gradient | Uses Monte Carlo sampled returns, applies likelihood ratio trick; pure policy gradient without value function               | High variance and sample inefficient; simpler but less stable than actor-critic or PPO                                     | Fundamental policy gradient algorithm, baseline for many RL studies |

# REINFORCE

REINFORCE is a method that is very close to the Policy Gradient Theorem formulation

- uses Advantage rather than the episode Return
- the Advantage subtracts a "baseline" from the episode Return

| Method    | Intermediate Reward Advantage $A_\tau$ ,<br>$\backslash\text{advseq}\_t = G_t - b_t$ | Single Trajectory Reward Advantage $A_\tau$<br>$\backslash\text{advseq} = R(\tau) - b$ |
|-----------|--|--|
| REINFORCE |  |  |

The purpose of subtracting a baseline is to

- reduce the magnitude of the gradients
- reduces the variance of the gradients
- and hence: the change in policy parameters  $\theta$

In single, end-of-episode reward, the Advantage is typically computed as

$$\text{advseq}_{\tau} = (\text{rewseq}(\tau) - b)$$

- $b$  is often the moving average (over trajectories  $\tau$ ) of  $\text{rewseq}(\tau)$
- advantage is the same for every step

With intermediate rewards, the Advantage is typically computed as

$$\text{advseq}_{\{\tau, t\}} = \left( G_{\{\tau, t\}} - b_{t\tau} \right)$$

- where  $b_{t\tau}$  is often a proxy for the *expected* Value function  
 $\text{statevalfun}_{\pi}(\text{stateseq})$  of state  $\text{stateseq}_{t\tau}$

It uses a Monte-Carlo method to estimate the gradient

- based on a sample of trajectories

## Pseudo code for REINFORCE

```
# REINFORCE training for LLM
for prompt in training_prompts:
    output = llm.generate(prompt)
    reward = evaluate_output(output) # Human or automated score
    logprob = llm.logprob(output, prompt)

    # Monte Carlo policy gradient update (no critic)
    baseline = compute_baseline() # Optional: running mean for variance reduction
    loss = -logprob * (reward - baseline)
    loss.backward()
    optimizer.step()
```

## **Discussion**

REINFORCE is very close to the vanilla Policy Gradient.

But it is considered high variance.

Subsequent methods will take explicit steps to reduce variance.

Notice that in the *mathematical* formulation of REINFORCE

$$G_{,\tau}$$

which is part of the Advantage

- is computed for *each trajectory*  $\tau$  independently
- the stochastic Transition Probability for a given State/Action pair (`\state, \act`)  
\$\$

**\transp({ \state', \rew | \state, \act })**

`\transp({ \stateseq{\tt+1}, \rewseq{\tt+1} | \state = \stateseq{\tt}, \act = \actseq{\tt} })`  
\$\$

is a *single sample* from a probability distribution for the State/Action pair (`\state, \act`).

This leads to *high variance* estimates of  $G_{,\tau}$

### Note

we observe the *effect* of the probability, not its value directly because we are model-free

In the *practical* implementation of REINFORCE

- we estimate  $G$
- over a *batch* of episodes  $\tau_1, \dots,$

so we have *multiple samples* from the probability distribution for the State/Action pair (`\state, \act`).

- hopefully leading to a lower variance approximation

# PPO

## Intuition

PPO optimizes a **clipped surrogate** loss

- which can be interpreted as variant of the standard policy gradient theorem.

The goal of the surrogate objective is to control how rapidly the policy changes from epoch to epoch of training.

- avoid large policy shifts
- variance reduction
- monotonic improvement of policy

# Surrogate Loss for PPO

The Surrogate Loss for PPO is

$$J_{\text{PPO}}(\theta) = \mathbb{E} \left[ \min \left( r_t(\theta) \hat{A}_t; \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{\text{advseq}}_t \right) \right]$$

We will explain how this potentially intimidating equation came about.

We express the

- relative change in policy (for a given State/Action pair)
- between epochs
- with the *probability ratio*

$$r(\theta) = \frac{\pi_\theta(\text{actseq} | \text{stateseq})}{\pi_{\theta_{\text{old}}}(\text{actseq} | \text{stateseq})}$$

where  $\pi_{\theta_{\text{old}}}$  is the *reference policy*

The reference policy is the policy in effect at the start of each epoch of training

- before this epoch modifies the policy parameters

During an epoch of training

- trajectories are generated using the reference policy
- the policy parameters are updated based on the surrogate loss for these trajectories
- the updated policy becomes the reference policy for the next epoch

By keeping the probability ratio close to 1, we constrain the Gradient Ascent update step to a small change in policy.

How do we keep the probability ratio close to 1 ?

- by using clipping to constrain it to the range  $[1 - \epsilon, 1 + \epsilon]$  for small  $\epsilon$

this is the *clipped* part of the clipped surrogate objective.

Formally

The *surrogate loss* in PPO is:

---

$$\begin{aligned} L_{\text{sur}}(\theta) &= \mathbb{E}_{\pi_\theta} \left[ r(\theta) \hat{\text{advseq}}_t \right] \\ &= \mathbb{E}_{\pi_\theta} \left[ \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \right] \end{aligned}$$

where

- $r(\theta)$  is the probability ratio,
  - $\hat{\text{advseq}}_t$  is the advantage estimate at step  $t$ .
-

We modify the Surrogate Loss to limit the effect of the Probability Ratio.

The *clipped surrogate loss* is

$$L_{\text{clip}}(\theta) = \mathbb{E}_{r \sim p_\theta} [\min\{\max(1 - \hat{p}_\theta(r), 0), \hat{p}_\theta(r)\}]$$

By minimizing the clipped surrogate loss

- we maximize  $J(\theta)$  (our goal in Gradient Ascent)
  - in a controlled manner
-

## Relation to the Unified Gradient Formulation

The (unclipped) Surrogate Loss does not resemble the Unified Gradient Formulation.

- absence of the term

$$\log \pi_\theta(\text{actseq}_\tau | \text{stateseq}_{\tau,t})$$

that multiplies the Advantage.

Moreover, there is an extra term

$$r(\theta) = \frac{\pi_\theta(\text{actseq} | \text{stateseq})}{\pi_{\theta_{\text{old}}}(\text{actseq} | \text{stateseq})}$$

in the form of the Probability Ratio.

But, these terms appear inside the Expectation of the Gradient and, by using the Likelihood Ratio trick

$$\begin{aligned}\nabla_{\theta} r(\theta) &= r(\theta) \nabla_{\theta} \log r(\theta) \\ &= r(\theta) \nabla_{\theta} (\log \pi_{\theta}(\text{\textbackslash actseq}_1 | \text{\textbackslash stateseq}) - \log \pi_{\text{old}}(\text{\textbackslash actseq}_1 | \text{\textbackslash stateseq})) \\ &= r(\theta) \nabla_{\theta} \log \pi_{\theta}(\text{\textbackslash actseq}_1 | \text{\textbackslash stateseq})\end{aligned}$$

Thus, taking the gradient of the Probability Ratio

- yields the usual  
 $\log \pi_{\theta}(\text{\textbackslash actseq}_{\tau, t} | \text{\textbackslash stateseq}_{\tau, t})$

term.

---

We still have an extra  $r_c(\theta)$  multiplicative term (relative to the Unified Gradient Formulation).

So we don't obtain an *identical* expression

- but we note that, with clipping,  $r_c(\theta) \approx 1$

## Advantage for PPO

To complete the presentation in terms of the Unified Gradient formulation, we define the Advantage.

| Reward Setting       | Advantage Estimate  |
|----------------------|---|
| Intermediate Rewards | $\hat{\text{advseq}}_t = G_t \text{ or GAE}$<br>$- \text{statevalfun}_\pi(\text{stateseq})$ |
| Single Total Reward  | $\hat{\text{advseq}}_t = R(\tau)$<br>$- b$  |

## Common ways to compute Advantage Estimate $\hat{\text{advseq}}_t$ for Intermediate Rewards case

- Using Return-to-go and value estimate:

$$\hat{\text{advseq}}_t = G_t - V(s_t)$$

where  $G_t$  is the discounted sum of rewards (return-to-go) starting from time step  $t$ , and  $V(s_t)$  is the estimated state value.

---

- Generalized Advantage Estimation (GAE):

$$\hat{\text{advseq}}_t = \sum_{l=0}^{\infty} (\gamma \lambda)^l \delta_{t+l}$$

where

$$\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$$

with

- $\gamma$  as the discount factor,
- $\lambda \in [0, 1]$  controlling the bias-variance trade-off.

- When  $\lambda = 0$ , GAE reduces to the **1-step Temporal Difference (TD) advantage**:
- When  $\lambda = 1$ , GAE reduces to the **Monte Carlo advantage**:

GAE smoothly interpolates between high-bias low-variance and low-bias high-variance advantage estimates, making it very effective in practice.

## Pseudo code for PPO

```
# PPO training for LLM
for prompt in training_prompts:
    output = llm.generate(prompt)
    reward = evaluate_output(output)
    logprob_old = llm.logprob(output, prompt) # From previous policy
    value = critic(output, prompt) # Critic gives value baseline

    # Calculate advantage
    advantage = reward - value

    # Compute importance ratio
    logprob_new = llm_new.logprob(output, prompt)
    ratio = exp(logprob_new - logprob_old)

    # Clipped surrogate objective for stability
    clip_epsilon = 0.2
    loss1 = ratio * advantage
    loss2 = clip(ratio, 1-clip_epsilon, 1+clip_epsilon) * advantage
    loss = -min(loss1, loss2)
    loss.backward()
```

## PPO with batches

```
for epoch in range(num_epochs):
    for batch in data_loader:
        # Generate trajectories
        contexts = batch["contexts"]
        responses = model.sample(contexts) # generate outputs

        # Score trajectories with reward model
        rewards = reward_model.score_batch(contexts, responses)

        # Compute value estimates and advantages (using GAE or similar)
        values = value_model.predict(contexts, responses)
        advantages = compute_advantages(rewards, values)

        # Get old policy log-probs for PPO ratio calculation
        old_log_probs = model.log_prob(contexts, responses).detach()

        # Forward pass to get new log probabilities
        new_log_probs = model.log_prob(contexts, responses)

        # Calculate PPO ratio
```

Note that all the mathematical operations

- are performed in parallel on each example in the batch
- the per-example loss
  - is reduced to a single scalar
  - via the `.mean()`

```
loss = -torch.min(surrogate1, surrogate2).mean()
```

# Discussion

PPO takes explicit steps

- clipped probability ratio

to reduce variance of Gradient estimates.

It is robust and is a very common method when performing RL Tuning.

## Trust Region Policy Optimization (TRPO) Theory (Brief)

We stated that, constraining the probability ratio

$$r_t(\theta) = \frac{\pi_\theta(\text{actseq} | \text{stateseq})}{\pi_{\theta_{\text{old}}}(\text{actseq} | \text{stateseq})}$$

to the small range  $[1 - \epsilon, 1 + \epsilon]$

is the key to the monotonic improvement of the optimization objective.

This is based on *Trust Region Policy Optimizaton (TRPO)*.

We state TRPO briefly.

TRPO formulates policy optimization as a constrained optimization problem to ensure stable and monotonic policy improvement.

It maximizes a surrogate objective subject to a constraint on how much the new policy can deviate from the old policy:

$$\max_{\theta} \quad \mathbb{E}_{s,a \sim \pi_{\theta_{\text{old}}}} \left[ \frac{\pi_{\theta}(a|s)}{\pi_{\theta_{\text{old}}}(a|s)} A^{\pi_{\theta_{\text{old}}}}(s, a) \right]$$

subject to the trust region constraint:

$$\mathbb{E}_{s \sim \pi_{\theta_{\text{old}}}} [D_{\text{KL}}(\pi_{\theta_{\text{old}}}(\cdot|s) \| \pi_{\theta}(\cdot|s))] \leq \delta$$

where:

- $\pi_{\theta_{\text{old}}}$  is the old (reference) policy,
- $\pi_{\theta}$  is the new policy parameterized by  $\theta$ ,
- $A^{\pi_{\theta_{\text{old}}}}(s, a)$  is the advantage function with respect to the old policy,
- $D_{\text{KL}}(\cdot \| \cdot)$  is the Kullback-Leibler (KL) divergence measuring the difference between the old and new policies,
- $\delta$  is a small positive constant controlling the maximum allowed policy update.

The trust region constraint (the KL divergence term)

- **limits the size of the policy update**
- to ensure the new policy does not differ drastically from the old policy
- preventing performance collapse.

This facilitates **monotonic improvement** in policy performance.

PPO can be interpreted as a practical approximation of TRPO

- that replaces the hard KL constraint
- with a clipped probability ratio

in the objective.

---

# **GRPO**

GRPO is a Policy Gradient based methods that uses *preferences* rather than Rewards.

It will be presented in a separate module with other Preference based methods.

In [2]: `print("Done")`

Done

